

Modeling zero-inflated explanatory variables in hybrid Bayesian network classifiers for species occurrence prediction

Ana D. Maldonado, Pedro A. Aguilera, Antonio Salmerón

Published in:

Environmental Modelling & Software

DOI (link to publication from publisher):

<https://doi.org/10.1016/j.envsoft.2016.04.003>

Publication date:

2016

Document version:

Accepted author manuscript, peer reviewed version

Citation for published version:

Maldonado, A. D., Aguilera, P. A., & Salmerón, A. (2016). Modeling zero-inflated explanatory variables in hybrid Bayesian network classifiers for species occurrence prediction. *Environmental Modelling & Software*, 82, 31-43.
<https://doi.org/10.1016/j.envsoft.2016.04.003>

Modeling zero-inflated explanatory variables in hybrid Bayesian network classifiers for species occurrence prediction

A.D. Maldonado^{a,*}, P.A. Aguilera^b, A. Salmerón^a

^a*Department of Mathematics, University of Almería, Almería, Spain*

^b*Informatics and Environment Laboratory, Department of Biology and Geology, University of Almería, Almería, Spain*

Abstract

Datasets with an excessive number of zeros are fairly common in several disciplines. The aim of this paper is to improve the predictive power of hybrid Bayesian network classifiers when some of the explanatory variables show a high concentration of values at zero. We develop a new hybrid Bayesian network classifier called *zero-inflated tree augmented naive Bayes* (Zi-TAN) and compare it with the already known *tree augmented naive bayes* (TAN) model. The comparison is carried out through a case study involving the prediction of the probability of presence of two species, the fire salamander (*Salamandra salamandra*) and the Spanish Imperial Eagle (*Aquila adalberti*), in Andalusia, Spain. The experimental results suggest that modeling the explanatory variables containing many zeros following our proposal boosts the performance of the classifier, as far as species distribution modeling is concerned.

Keywords: hybrid Bayesian networks, mixtures of truncated exponentials, zero excess treatment, species distribution modeling

*Corresponding author

Email addresses: `ana.d.maldonado@ual.es` (A.D. Maldonado), `aguilera@ual.es` (P.A. Aguilera), `antonio.salmeron@ual.es` (A. Salmerón)

Software availability

The algorithms introduced in this paper have been implemented within the Elvira environment for probabilistic graphical models (Elvira Consortium, 2002), which is a free open source software programmed in Java. The software, including the necessary scripts for replicating the experiments reported in this paper, can be downloaded from the website

<http://www.ual.es/personal/amg457/downloads>

where the datasets used in the paper are also available for download. Both the software and the data are contained in a single zip file, which includes README files with the necessary instructions. The size of the zip file is 6.5 MB.

The software has been compiled with Oracle Java™ SE version 1.8.0.45 build 14. It is ready to work in Windows, Mac and Linux platforms. The datasets provided are in `dbc` format, which is a plain text format used by the Elvira software, which provides facilities for exporting it to `csv` format.

1. Introduction

Environmental datasets tend to present a number of problems, which must be detected and solved in order to obtain plausible results (Ancelet et al., 2010; Lecomte et al., 2013). One of these problems is the presence of data with highly skewed frequency distributions containing an excessive number of zeros. As a consequence, the data do not follow a standard distribution and the application of the usual analysis techniques may yield inaccurate parameter estimates and misleading inferences (Martin et al., 2005). Examples of data with many zeros often occur in different fields, including environmental sciences (Potts and Elith, 2006; Kamarianakis et al., 2008; Dorevitch et al., 2011), ecology (Damgaard, 2008; Wenger and Freeman, 2008; Calama et al., 2011), epidemiology (Böhning et al., 1999; Ngatchou-Wandji and Paris, 2011), genetics (Varona and Sorensen, 2010), biochemistry (Nie et al., 2006; McDavid et al., 2013) or economy (Edmeades and Smale, 2006; Solé-Auró et al., 2012).

The algorithms developed to deal with zero excess are typically focused on the dependent variable. The most popular models are usually extensions of the *Generalized Linear Models*, comprising *zero-inflated Binomial* (ZIB) model (Hall, 2000) for binary variables; *zero-inflated Poisson* (ZIP) (Lambert, 1992), *zero-inflated Negative Binomial* (ZINB) (Greene, 1994), *Poisson*

21 *hurdle* and *Negative Binomial hurdle* (Cragg, 1971; Mullahy, 1986) for dis-
22 crete variables; and *delta models* or *compound Poisson process* for continuous
23 variables (Ancelet et al., 2010; Lecomte et al., 2013). Generalizing, these
24 models are combinations of probability distributions which separately model
25 the occurrence of zeros and the rest of the domain of the variable of interest.
26 These models are appropriate for handling response variables with high con-
27 centration of zeros and have been shown to outperform methodologies that
28 assume the dependent variable to have a standard distribution (Martin et al.,
29 2005).

30 However, the distribution of the explanatory variables has not typically
31 been of concern and therefore methodologies for dealing with explanatory
32 variables containing high concentration of zeros have not been studied so far.
33 Notwithstanding, accurately modeling the distribution of the explanatory
34 variables is crucial in models such as Bayesian network classifiers, that have
35 been successfully utilized in species distribution analysis (Aguilera et al.,
36 2010). Unfortunately, the methods described above for handling zero excess
37 are not directly applicable to Bayesian network classifiers because they are
38 not designed for modeling conditional distributions and in the case of con-
39 tinuous variables, they rely on distributions that are not compatible with
40 Bayesian network algorithms, as is the case of the Gamma distribution.

41 Bayesian networks (BNs) belong to the so-called *probabilistic graphical*
42 *models* and roughly speaking they are compact representations of joint prob-
43 ability distribution over a set of variables whose independence relations are
44 encoded by the structure of an underlying directed acyclic graph (Pearl,
45 1988). When a BN hosts discrete and continuous variables simultaneously,
46 it is called a *hybrid* BN. However, not every kind of distribution is compati-
47 ble with the factorization encoded in a hybrid Bayesian network. One of the
48 most flexible models is based on the use of *mixtures of truncated exponentials*
49 (MTEs), introduced by Moral et al. (2001), generalized later by Shenoy and
50 West (2011) and Langseth et al. (2012).

51 A hybrid BN classifier is just a BN where one of the variables is the
52 *class* (which is discrete) while the others (discrete or continuous) are the
53 explanatory variables, also called *features* (Aguilera et al., 2011). Typically,
54 when facing classification problems only restricted network structures are
55 considered, such as naive Bayes (NB) or tree augmented naive Bayes (TAN)
56 (Friedman et al., 1997). The NB model assumes that the explanatory vari-
57 ables are independent of each other given the class variable, while the TAN
58 model relaxes that assumption by allowing some dependencies among the

59 features. Within the Environmental Sciences area, the NB model appears
60 to be more popular (Markus et al., 2010; Aguilera et al., 2013; Fytilis and
61 Rizzo, 2013; Roperio et al., 2014, 2015) than the TAN model (Aguilera et al.,
62 2010; Maldonado et al., 2015).

63 In this paper we address the problem of having a high concentration of
64 zeros in explanatory variables of hybrid BN classifiers. More precisely, we
65 introduce a new model called *zero-inflated* TAN (Zi-TAN) that extends the
66 hybrid BN classifier proposed by Aguilera et al. (2010) by explicitly modeling
67 the zero values. We show how the new model outperforms the formerly used
68 hybrid BN classifier in a case study related to Species Distribution Models
69 (SDM). In the case study, environmental variables are used as explanatory
70 variables of the species occurrence, including climate, land use, soil and lithol-
71 ogy. Depending on the scale, these variables may contain a large proportion
72 of zeros which justifies the development of the new model.

73 The remainder of the paper is organized as follows. We describe Bayesian
74 network classifiers and our baseline model, the TAN, in Section 2. Section 3 is
75 devoted to the methodological aspects of our new proposal. The performance
76 of the new model is analyzed in a case study involving two species in Section 4.
77 The paper ends with conclusions in Section 5.

78 2. Bayesian networks for classification

79 A Bayesian network (BN) is a statistical multivariate model for a set of
80 variables $\mathbf{X} = \{X_1, \dots, X_n\}$, which is defined in terms of two components:

- 81 • Qualitative component: A directed acyclic graph (DAG) where each
82 vertex represents one of the variables in the model, and so that the
83 presence of an edge linking two variables indicates the existence of
84 statistical dependence between them.
- 85 • Quantitative component: A conditional distribution $p(x_i|pa(x_i))$ for
86 each variable X_i , $i = 1, \dots, n$ given its parents in the graph, denoted
87 as $pa(X_i)$.

88 The joint distribution of the variables in the network is therefore repre-
89 sented in a factorized way as

$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i|pa(x_i)) \quad \forall x_1, \dots, x_n \in \Omega_{X_1, \dots, X_n} \quad (1)$$

90 where Ω_{X_i} represents the set of all possible values of variable x_i and $pa(x_i)$
 91 denotes an instantiation of the parents of X_i .

92 Hybrid BNs can handle both discrete and continuous data without impos-
 93 ing restrictions on the interactions among the variables thanks to the devel-
 94 opment of models such as the *Mixtures of Truncated Exponentials* (MTEs)
 95 developed by Moral et al. (2001). The MTE model is characterized by a
 96 function defined as follows.

97 **Definition 1.** (MTE potential) *Let \mathbf{X} be a mixed n -dimensional random*
 98 *vector. Let $\mathbf{W} = (W_1, \dots, W_d)$ and $\mathbf{Z} = (Z_1, \dots, Z_c)$ be the discrete and*
 99 *continuous parts of \mathbf{X} , respectively, with $c + d = n$. We say that a function*
 100 *$f : \Omega_{\mathbf{X}} \mapsto \mathbb{R}_0^+$ is a Mixture of Truncated Exponentials potential (MTE*
 101 *potential) if for each fixed value $\mathbf{w} \in \Omega_{\mathbf{W}}$ of the discrete variables \mathbf{W} , the*
 102 *potential over the continuous variables \mathbf{Z} is defined as:*

$$f(\mathbf{z}) = a_0 + \sum_{i=1}^m a_i \exp \left\{ \sum_{j=1}^c b_i^{(j)} z_j \right\} \quad (2)$$

103 for all $\mathbf{z} \in \Omega_{\mathbf{Z}}$, where $a_i, i = 0, \dots, m$ and $b_i^{(j)}, i = 1, \dots, m, j = 1, \dots, c$ are
 104 real numbers. We also say that f is an MTE potential if there is a partition
 105 D_1, \dots, D_k of $\Omega_{\mathbf{Z}}$ into hypercubes and in each D_i , f is defined as in Eq. (2).

106 An MTE function is an *MTE density* if it integrates to 1. A *conditional*
 107 *MTE density* can be specified by dividing the domain of the conditioning
 108 variables and giving an MTE density of the conditioned variable for each
 109 configuration of splits of the other variables. The more the intervals used to
 110 divide the domain of the continuous variables, the better the MTE model
 111 accuracy but in exchange of a higher number of parameters. To estimate the
 112 parameters of MTE densities, we followed the approach recently introduced
 113 by Langseth et al. (2014), which is based on least squares optimization, but
 114 limiting the number of exponential terms to 2, i.e., $m = 2$ in Eq. (2), in order
 115 to keep the complexity of the models moderate.

116 Hybrid BNs can also be modeled by discretizing the continuous variables,
 117 so that all the existing methodology for discrete BNs can be applied with
 118 no further modification. The most prominent proposal in this direction is
 119 the so-called *dynamic discretization* (Neil et al., 2007) which seeks for better
 120 representations of high density areas throughout the inference process. The
 121 problem with discretization is to balance the desire for high accuracy in the
 122 approximations with a reasonable complexity of the resulting models. A

123 study of the complexity of the MTE approach versus discretization can be
 124 found in (Rumí and Salmerón, 2007; Langseth et al., 2009).

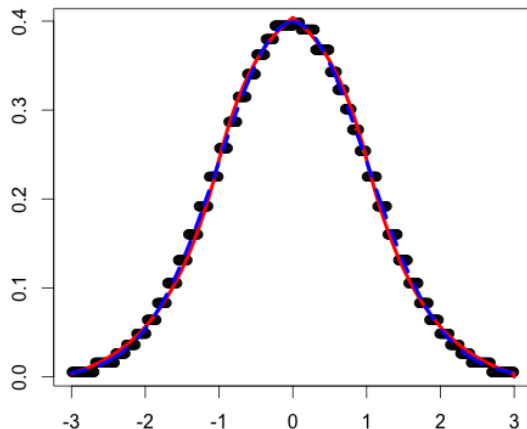


Figure 1: Plot of a standard normal density (dashed blue line) overlaid on an approximation using MTEs (solid red line) and dynamic discretization (solid black line).

125 As an illustration of the potential advantages of the MTE approach versus
 126 dynamic discretization, consider the problem of approximating a standard
 127 normal density using both approaches. An approximation using MTEs is
 128 given by Cobb et al. (2006) as

$$f(x) = \begin{cases} -0.017203 + 0.9309604e^{1.27x} & \text{if } -3 \leq x < -1, \\ 0.442208 - 0.038452e^{-1.64x} & \text{if } -1 \leq x < 0, \\ 0.442208 - 0.038452e^{1.64x} & \text{if } 0 \leq x < 1, \\ -0.017203 + 0.9309604e^{-1.27x} & \text{if } 1 \leq x < 3, \end{cases} \quad (3)$$

129 An approximation of the standard normal density using dynamic dis-
 130 cretization can be obtained using the AgenaRisk software.¹ Figure 1 shows
 131 both approximations overlaid on the plot of the standard normal density be-
 132 tween -3 and 3. The plot illustrates how using MTEs the approximation is

¹<http://www.agenarisk.com>

133 smooth while the discretized version is a staircase function. Hence, it is pos-
 134 sible to obtain more accurate approximations using fewer parameters with
 135 MTEs in general. In this case, the MTE approximation in Equation (3) has
 136 12 parameters while the discretized approximation provided by AgenaRisk
 137 has 50 parameters. We have also computed the mean absolute error of both
 138 approximations, obtaining a value of 0.0045 for the MTE model and 0.0055
 139 for the discretized one.

140 Hence, the potential benefits of using MTEs instead of discretized models
 141 are: (i) they provide, in general, more accurate approximations using fewer
 142 parameters, which leads to more compact models that require fewer param-
 143 eters to be estimated from data, (ii) they can easily represent variables whose
 144 nature is not discrete nor continuous, as we will discuss in Section 3. Fur-
 145 thermore, a discretized model can be seen as a particular case of an MTE
 146 where only parameter a_0 in Equation (2) is different from 0.

147 A Bayesian network can be used as a classifier if it contains a class variable
 148 C and a set of continuous or discrete explanatory variables X_1, \dots, X_n , where
 149 an object with observed features x_1, \dots, x_n will be classified as belonging to
 150 class $c^* \in \Omega_C$ obtained as

$$c^* = \arg \max_{c \in \Omega_C} p(c|x_1, \dots, x_n),$$

151 where Ω_C denotes the set of all possible values of C .

152 Considering that $p(c|x_1, \dots, x_n)$ is proportional to $p(c) \times p(x_1, \dots, x_n|c)$,
 153 the specification of an n dimensional distribution for X_1, \dots, X_n given C is
 154 required in order to solve the classification problem, which implies a consid-
 155 erable computational cost, as the number of parameters necessary to specify
 156 a joint distribution is exponential in the number of variables, in the worst
 157 case. However, this problem is simplified if we take advantage of the factor-
 158 ization encoded by the BN. Since building a network without any structural
 159 restriction is not always feasible (they might be as complex as the above
 160 mentioned joint distribution), networks with fixed or restricted and simple
 161 structures are utilized instead when facing classification tasks. The extreme
 162 case is the naive Bayes (NB) structure, where all the feature variables are
 163 considered independent given C , as depicted in Fig. 2(a). The strong as-
 164 sumption of independence behind NB models is somewhat compensated by
 165 the reduction in the number of parameters to be estimated from data, since
 166 in this case, it holds that

$$p(c|x_1, \dots, x_n) \propto p(c) \prod_{i=1}^n p(x_i|c) , \quad (4)$$

167 which means that, instead of one n -dimensional conditional density, n one-
 168 dimensional conditional densities must be estimated.

169 In TAN models, more dependencies are allowed, expanding the NB struc-
 170 ture by permitting each feature to have one more parent besides C . It is
 171 illustrated in Fig. 2(b). The increase in complexity, in both the structure
 172 and the number of parameters, results in richer and more accurate models in
 173 general (Friedman et al., 1997).

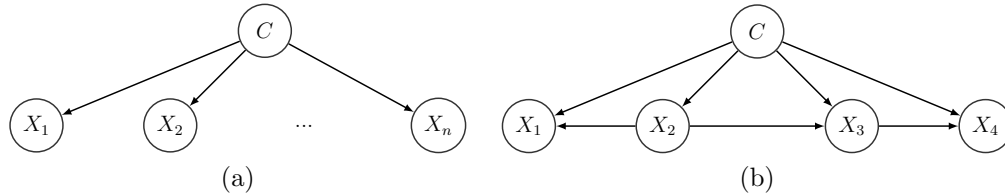


Figure 2: Structure of naive Bayes (a) and TAN (b) classifiers.

174 In general, there are several possible TAN structures for a given set of
 175 variables. The way to choose among them is to construct a maximum weight
 176 spanning tree containing the features, where the weight of each edge is the
 177 mutual information between the linked variables, conditional on the class
 178 (Friedman et al., 1997; Fernández et al., 2007). The mutual information
 179 between features X_i and X_j given the class is defined as

$$I(X_i, X_j|C) = \sum_{x_i, x_j, c} \log \frac{p(x_i, x_j|c)}{p(x_i|c)p(x_j|c)} . \quad (5)$$

180 The details for constructing a TAN classifier model are given in Algorithm 1.

181 3. Zero-inflated TAN based on mixtures of truncated exponentials

182 In environmental datasets, it is common to find variables with a high con-
 183 centration of observations at a single repeated value. This makes the mod-
 184 eling of the probability distribution for such variables a problematic task.
 185 As an example, consider the histogram on the left panel of Fig. 3. It repre-
 186 sents the distribution of *eutric regosols*, used in the case study in Section 4,

Algorithm 1: TAN classifier

Input: A dataset D with variables X_1, \dots, X_n, C .

Output: A TAN classifier with root variable C and features X_1, \dots, X_n .

- 1 Calculate the conditional mutual information $I(x_i, x_j|c)$ between each pair of attributes, $i \neq j$.
 - 2 Construct a complete undirected graph with nodes X_1, \dots, X_n and label each link connecting X_i to X_j by $I(x_i, x_j|c)$.
 - 3 Build a maximum weighted spanning tree \mathcal{T} .
 - 4 Transform \mathcal{T} into a directed tree by choosing a root variable, C , and setting the direction of every link to be outward from it.
 - 5 Construct a new network \mathcal{G} , with node C being connected to each X_i and nodes X_1, \dots, X_n having the same links as in \mathcal{T} .
 - 6 Estimate an MTE density for C , and a conditional MTE density for each X_i , $i = 1, \dots, n$ given its parents in \mathcal{G} .
 - 7 Let P be a set of estimated densities.
 - 8 Let TAN be a Bayesian network with structure \mathcal{G} and distribution P .
 - 9 **return** TAN.
-

187 including values equal to zero, which represent 667 out of 887 values. The
188 distribution is so concentrated at zero that the histogram provides no valu-
189 able information for values above zero. By excluding the zero values, the
190 resulting histogram is represented on the right panel of Fig. 3. It is apparent
191 that the distribution of the values greater than zero is far from being uni-
192 form, and therefore modeling it accurately can provide benefits in prediction
193 tasks.

194 Actually, the situation described above is somehow motivated by the fact
195 that the variable is not really discrete nor continuous. Instead, one can
196 consider that there is some probability mass allocated at 0, and the rest of
197 the probability mass is described by a density function. Formally, a variable
198 that is not discrete nor continuous is called a *mixed variable* in Statistics.
199 More precisely, a random variable is mixed if its distribution function has
200 discontinuity jumps at a countable number of points, and it is continuously
201 increasing at least in one interval of values of the variable.

202 As an example, let $g(x) \geq 0$ for $0 < x \leq 1$ be any non-negative real
203 function such that $\int_0^1 g(x)dx = 1 - p$, with $0 < p < 1$. Then, a random
204 variable X with density function

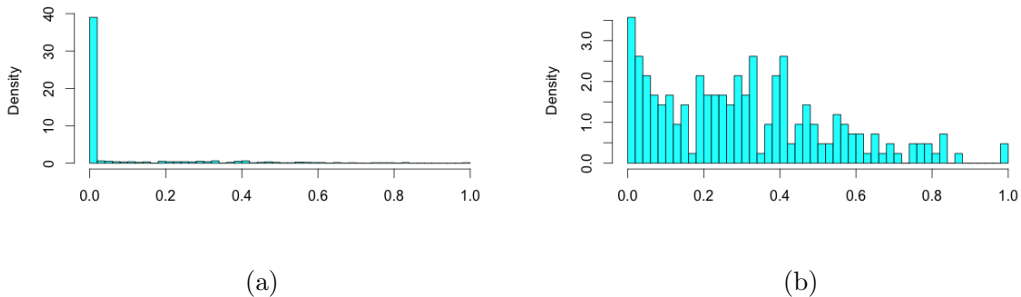


Figure 3: Histogram for the proportion of *eutric regosols* including (a) and excluding (b) values equal to zero.

$$f(x) = \begin{cases} p & \text{if } X = 0 \\ g(x) & \text{if } 0 < x \leq 1, \end{cases} \quad (6)$$

205 is a mixed random variable.

206 Considering the variable represented in Fig. 3, p would correspond to the
 207 fraction of observations allocated at 0 (i.e. the leftmost bar in the left panel),
 208 while $g(x)$ would correspond to the rest of the histogram or, equivalently, to
 209 the histogram on the right panel of the figure, which can be considered as the
 210 result of zooming in the initial histogram for the values of X strictly greater
 211 than 0. From now on, we will say that a mixed random variable whose density
 212 can be written as in Eq. (6) is a *zero-inflated random variable*. Note that we
 213 are considering, without loss of generality, that zero-inflated variables take
 214 values on $[0, 1]$. Variables with a different support can be re-scaled.

215 Zero-inflated random variables have not previously been considered in
 216 hybrid Bayesian network literature in general, nor in MTEs in particular.
 217 However, they can be easily accommodated within MTE models by incorporat-
 218 ing *artificial variables*. More precisely, our proposal consists in including
 219 an artificial variable X^* for each mixed variable X in the network, where X^*
 220 has no parents and X is its only child. The artificial variable is defined as
 221 follows:

$$X^* = \begin{cases} 0 & \text{if } X = 0 \\ 1 & \text{otherwise,} \end{cases} \quad (7)$$

222 and its probability function is

$$f(x^*) = P(X^* = x^*) = \begin{cases} p & \text{if } x^* = 0 \\ 1 - p & \text{if } x^* = 1, \end{cases} \quad (8)$$

223 where p is as in Eq. (6). Note that $f(x^*)$ is trivially an MTE, according to
224 Eq. (2).

225 The conditional distribution of X given X^* is

$$f(x|x^*) = \begin{cases} 1 & \text{if } x^* = 0, x = 0 \\ \frac{1}{1-p}g(x) & \text{if } x^* = 1, 0 < x \leq 1, \end{cases} \quad (9)$$

226 with p and $g(x)$ as in Eq. (6). Again, $f(x|x^*)$ is an MTE whenever $g(x)$
227 is an MTE as well. Note that so far we have made no assumptions about
228 $g(x)$ beyond those required for the corresponding density function being well
229 defined - see Eq. (6). We will see later in Definition 2 that $g(x)$ plays the
230 role of the conditional MTE distributions in a TAN classifier. The following
231 proposition states that the introduction of artificial variables does not modify
232 the marginal distribution of the zero-inflated variable.

233 **Proposition 1.** *Let X^* be a binary random variable with probability function*
234 *as in Eq. (8) and let X be a random variable whose distribution conditional*
235 *on X^* is as given in Eq. (9). Then, X is a zero-inflated random variable*
236 *with marginal distribution as in Eq. (6).*

Proof. The joint distribution of X and X^* is $f(x, x^*) = f(x|x^*)f(x^*)$, which
can be written as

$$f(x, x^*) = \begin{cases} p \times 1 & \text{if } x = 0, x^* = 0 \\ \frac{1}{1-p}g(x) \times (1-p) & \text{if } 0 < x \leq 1, x^* = 1 \end{cases}$$

which amounts to

$$f(x, x^*) = \begin{cases} p & \text{if } x = 0, x^* = 0 \\ g(x) & \text{if } 0 < x \leq 1, x^* = 1. \end{cases}$$

237 Therefore, the marginal distribution for X is obtained by marginalizing out
238 X^* as follows:

$$f(x) = \sum_{x^*=0}^1 f(x, x^*) = \begin{cases} p & \text{if } x = 0 \\ g(x) & \text{if } 0 < x \leq 1, \end{cases}$$

239 which matches Eq. (6). □

240 Our methodological proposal consists of including zero-inflated random
 241 variables in TAN classifiers, resulting in a new Bayesian network classifier
 242 formally defined as follows.

243 **Definition 2.** *Let \mathcal{T} be a TAN classifier over class variable C and features*
 244 *X_1, \dots, X_n . Let $X_i, i \in I \subset \{1, \dots, n\}$ be a set of zero-inflated random*
 245 *variables. A zero-inflated TAN (Zi-TAN) classifier \mathcal{T}^* is obtained from \mathcal{T}*
 246 *by:*

- 247 1. *Inserting, for each variable $X_i, i \in I$, an artificial variable X_i^* as in*
 248 *Eq. (7) and a link $X_i^* \rightarrow X_i$.*
- 249 2. *Attaching to each node $X_i^*, i \in I$ a distribution as in Eq. (8).*
- 250 3. *Attaching to each node $X_i, i \in I$ with parents $\{Y_1, \dots, Y_m\}$ in \mathcal{T} , and*
 251 *conditional distribution $f(x_i|y_1, \dots, y_m)$ in \mathcal{T} , a new conditional distri-*
 252 *bution*

$$f(x_i|x_i^*, y_1, \dots, y_m) = \begin{cases} 1 & \text{if } x_i^* = 0, x_i = 0 \\ \frac{1}{1-p}g(x_i|y_1, \dots, y_m) & \text{if } x_i^* = 1, 0 < x_i \leq 1, \end{cases} \quad (10)$$

253 where p is the proportion of values of X_i equal to 0 and $g(x_i|y_1, \dots, y_m)$
 254 is a conditional MTE density for X_i given $\{Y_1, \dots, Y_m\}$ learnt from the
 255 same sample as $f(x_i|y_1, \dots, y_m)$ but excluding the values where $X_i = 0$.

256 Notice that the new conditional distributions defined in Eq. (10) are of
 257 class MTE as long as the distributions in the original TAN model are MTEs
 258 as well. Also, the role of $g(x_i|y_1, \dots, y_m)$ corresponds to that of $g(x)$ in
 259 Eq.(6). Such conditional distributions are learnt making use of the procedure
 260 introduced by Langseth et al. (2014).

261 **Example 1.** *Consider the TAN structure in Figure 2(b). Assume that X_1*
 262 *and X_4 are zero-inflated random variables. The corresponding Zi-TAN struc-*
 263 *ture, according to Definition 2, is shown in Figure 4.*

264 The insertion of the artificial variables when constructing the Zi-TAN
 265 means that the new model can be factorized as a sum of TAN models, one
 266 per each combination of values of the artificial variables. However, from a
 267 practical point of view it is not a problem, as the joint distribution over the

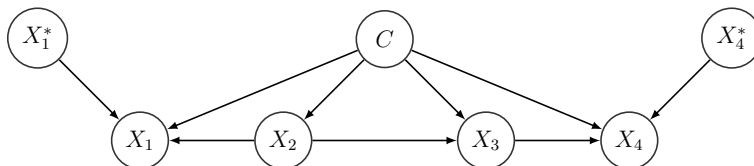


Figure 4: An example of a Zi-TAN classifier structure, obtained from Figure 2(b) assuming that X_1 and X_4 are zero-inflated random variables. X_1^* and X_4^* are their respective artificial variables.

268 class variables and the features is not affected, as shown in Proposition 2.
 269 Recall that the aim of the Zi-TAN model is not to modify the underlying
 270 distribution over the variables in the domain being analyzed, but rather to
 271 express it in a way that permits overcoming the problem of high concentration
 272 of values at zero.

273 **Proposition 2.** *Let \mathcal{T} be a TAN classifier over class variable C and fea-*
 274 *tures X_1, \dots, X_n , and \mathcal{T}^* be a Zi-TAN classifier constructed as in Defini-*
 275 *tion 2. Then, \mathcal{T}^* encodes the same probability distribution as \mathcal{T} over vari-*
 276 *ables $\{C, X_1, \dots, X_n\}$.*

277 *Proof.* According to Proposition 1, marginalizing out each artificial variable
 278 X_i^* in \mathcal{T}^* yields a conditional distribution for X_i exactly equal to the one it
 279 had in \mathcal{T} . Therefore, after removing all the artificial variables in \mathcal{T}^* , both
 280 models become the same. \square

281 The details on how to build a Zi-TAN classifier from data are given in
 282 Algorithm 2. It relies on Definition 2 and Algorithm 1.

283 4. Case study

284 In this section, the methodology explained above is applied to SDMs.
 285 More precisely, we considered two case studies involving the Fire Salamander
 286 and the Spanish Imperial Eagle.

287 4.1. Study area

288 The study area is Andalusia, a region in southern Spain which occupies
 289 an area of 87 000 km² and whose latitude and longitude is between 36°N -
 290 38°44'N and 3°50'W - 0°34'E. As far as elevation is concerned, the study area
 291 ranges from 0 to 3460 meters above the sea level. The main mountain ranges

Algorithm 2: Zi-TAN classifier

Input: A dataset D with variables X_1, \dots, X_n, C . A set of indices $I \subset \{1, \dots, n\}$ of zero-inflated variables.

Output: A Zi-TAN classifier with root variable C and features $\{X_1, \dots, X_n\} \cup \{X_i | i \in I\}$.

1 Build a TAN model, \mathcal{T} , from dataset D , for variables X_1, \dots, X_n, C using Algorithm 1.

2 $\mathcal{T}^* \leftarrow \mathcal{T}$.

3 **for** $i \in I$ **do**

4 Create a new binary variable

$$X_i^* = \begin{cases} 0 & \text{if } X_i = 0 \\ 1 & \text{otherwise .} \end{cases}$$

5 Add a new column to dataset D , corresponding to X_i^*

6 Insert a new link $X_i^* \rightarrow X_i$ in \mathcal{T}^* .

7 Let p be the proportion of values $X_i^* = 0$ in D .

8 Attach to X_i^* in \mathcal{T}^* the distribution

$$f(x_i^*) = \begin{cases} p & \text{if } x_i^* = 0 \\ 1 - p & \text{if } x_i^* = 1, \end{cases}$$

9 Let $\{Y_1, \dots, Y_m\}$ be the parents of X_i in \mathcal{T} .

10 Let $f(x_i | y_1, \dots, y_m)$ be the conditional distribution of X_i given its parents in \mathcal{T} .

11 Estimate a new density $g(x_i | y_1, \dots, y_m)$ from the same data used to learn $f(x_i | y_1, \dots, y_m)$, but excluding the elements in the sample where $X_i = 0$.

12 Attach to X_i in \mathcal{T}^* a new conditional distribution

$$f(x_i | x_i^*, y_1, \dots, y_m) = \begin{cases} 1 & \text{if } x_i^* = 0, x_i = 0 \\ \frac{1}{1-p} g(x_i | y_1, \dots, y_m) & \text{if } x_i^* = 1, 0 < x_i \leq 1, \end{cases}$$

13 **return** \mathcal{T}^* .

292 of Andalusia are the Sierra Morena mountain range (in the North) and the
293 Baetic systems (in the South), which are separated by the Baetic depression,
294 the lowest territory in Andalusia (Figure 5). The flattest areas correspond to
295 the littoral and the Baetic depression, through which the Guadalquivir river
296 runs, and the steepest ones to the Baetic Systems, comprising the Prebaetic,
297 Subbaetic and Pennibaetic systems.

298 The geographic location determines Andalusia's climate, which belongs
299 to the Mediterranean domain. The Mediterranean climate alternates mild,
300 rainy and humid winters with dry and warm summers. The average annual
301 temperature usually does not drop below 15°C, as a consequence of the ocean
302 influence. On the other hand, precipitation (P) shows a high spatial vari-
303 ability, ranging from 170 mm/year to 2180 mm/year. In addition, potential
304 evapotranspiration (PET) ranges from about 300 mm in the eastern Baetic
305 systems to more than 1000 mm/year in both the Guadalquivir river area and
306 the eastern coast. As a result, the quotient of precipitation divided by PET,
307 i.e. the humidity index, varies along the study area from surplus or areas
308 with hydric excess ($P > PET$) to deficit or areas lacking in water resources
309 ($P < PET$).

310 Regarding land use, half of the study area corresponds to natural vegeta-
311 tion, followed by croplands (44%), urbanized areas (3%) and bodies of water
312 (3%). With reference to soil classification, cambisols are the most common
313 soil group (41%), followed by regosols (19%), vertisols (9%), litosols (8.2%),
314 luvisols (8.15%), fluvisols (5.5%), planosols (2.5%), xerosols (2%), solonchaks
315 (2%), arenosols (1%), dunes (0.12%) and histosols (0.01%). On the subject
316 of lithology, 64% of the Andalusia's crust consists of sedimentary rocks, fol-
317 lowed by metamorphic rocks (26%), plutonic rocks (6%) and volcanic rocks
318 (4%).

319 *4.2. Data description*

320 Data from different thematic maps (Table 1) were incorporated into a
321 geographic information system - ArcGis (ESRI® ArcMap™10.2.2). A 10
322 x 10 km grid with presence records of different species was superimposed
323 on the thematic maps in order to build a matrix composed of a number of
324 explanatory variables and 1 target variable. The coordinate system for all
325 the datasets is based on the European Terrestrial Reference System 1989
326 (ETRS89).

327 The percentage of each land use, soil and lithology variable within each
328 cell was calculated by dividing the area of each variable by the total cell

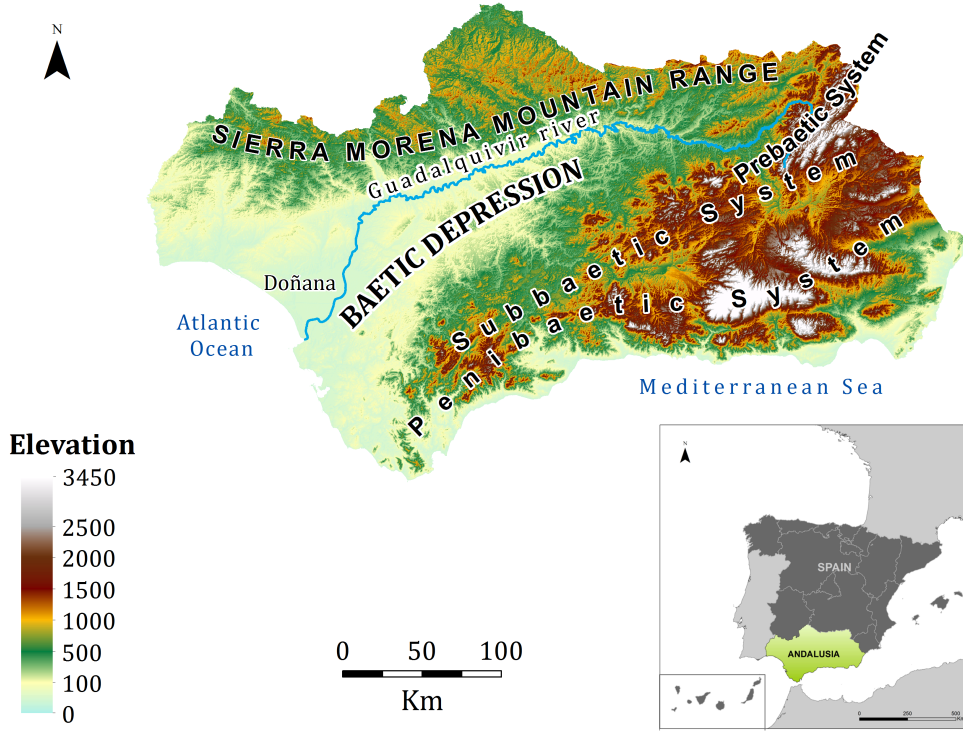


Figure 5: Study area, Andalusia, Spain (37°23'00"N 5°59'00"W).

329 area. In addition, the average of annual mean temperature, annual rainfall,
 330 annual PET and annual humidity index were calculated for the 30 year period
 331 1971-2000 for each cell in the grid. Moreover, the average elevation, slope
 332 and aspect of each cell was calculated from the Andalusian Digital Terrain
 333 Model (DTM). The average elevation and slope in each cell were determined
 334 by means of their arithmetic means. Since aspect is measured clockwise from
 335 0 to 360 degrees from the north, describing a full circle, its average (\bar{A}) was
 336 calculated as in Davis (1986)

$$\bar{A} = \arctan \left(\frac{\sum_{i=1}^n \sin \theta_i}{\sum_{i=1}^n \cos \theta_i} \right), \quad (11)$$

337 where θ_i is the aspect angle in each DTM pixel and n is the number of pixels.
 338 Once the value of each variable was calculated for the 10 x 10 km grid,

339 those cells occupied by less than 50% of terrestrial surface were removed.
340 Afterwards, a matrix composed of 156 variables taking values over 887 cells
341 was obtained, where the class variable is the presence of a particular species
342 and the remaining 155 variables are the explanatory variables (Table 1).
343 Finally, the data were rescaled to interval [0,1] in order to prevent numerical
344 instability problems.

345 We have considered an explanatory variable to be zero-inflated if more
346 than 50% of its observations are zeros. In this regard, 128 explanatory vari-
347 ables out of 155 came out to be zero-inflated. In order to test the *zero-inflated*
348 TAN (Zi-TAN) classifier, 2 species were chosen: 1) the Fire Salamander
349 (*Salamandra salamandra*; from now on referred to as salamander), which is
350 present in 300 out of 887 cells; therefore it was neither too present nor too
351 absent; and 2) the Spanish Imperial Eagle (*Aquila adalberti*; from now on
352 referred to as eagle), whose presence in the study area is imbalanced and
353 scarce, occurring in 54 out of 887 cells. Therefore, the 2 species are different
354 from the point of view of occurrence.

355 4.3. Variable selection

356 An excessive number of variables may decrease both the generalization
357 of the model by overfitting and its performance by introducing noise. Since
358 the dataset contained a high number of variables, a variable selection process
359 was carried out. For both species, the selection process was entirely carried
360 out by experts.

361 In the case of salamander, since the species prefers dark moist areas, with
362 rainfall being abundant, and is widely found in the Mediterranean forest at
363 medium-high elevations (Degani and Warburg, 1978), the selected variables
364 were rainfall, humidity index, dense oak woodlands, oak woodlands with
365 shrub, oak woodlands with herbaceous crops, woodlands with herbaceous
366 crops, grassland, olive groves, eutric regosols, calcaric regosols, eutric cam-
367 bisols, sand-silt-clay-gravel lithology type, slate-greywacke-sandstone lithol-
368 ogy type and volcano-sedimentary complex lithology type. Within this se-
369 lection, 7 out of 14 were zero-inflated explanatory variables.

370 Regarding eagle, the 3 main populations in Andalusia differ in terms
371 of environmental conditions, ranging from the mediterranean forest to low-
372 lying marshes near the ocean (González et al., 2008). Therefore, the selected
373 variables were temperature, rainfall, evapotranspiration, rainfed herbaceous
374 crops, oak woodlands with shrub, oak woodlands with herbaceous crops,

375 marshes, eutric regosols, albic arenosols, solonchaks, eutric cambisols, slate-
376 shale-greywacke-quartzite lithology type, slate-greywacke-sandstone lithol-
377 ogy type, sand lithology type and silt-clay lithology type. Within this selec-
378 tion, 10 out of 15 were zero-inflated explanatory variables.

379 4.4. Learning species distribution models

380 Once the variables were selected, we used Algorithm 1 for learning the
381 TAN model and Algorithm 2 for learning the Zi-TAN model.

382 A powerful advantage of Bayesian networks is their capability of predict-
383 ing a full specification of the posterior distribution of the class variable. In
384 the case of species distribution modeling, we are interested in determining
385 the probability of occurrence of a particular species, rather than giving a
386 fixed binary prediction.

387 The classifiers described in Sections 2 and 3 can be used to depict the
388 potential distribution area, defined as the *probability of presence of the species*
389 in each 10x10 km cell given the explanatory variables (X_1, \dots, X_n) included
390 in the model. Each cell in the map can be colored across a gradual white to
391 red color ramp, where low probability of presence (0-0.2) is represented with
392 the white color and high probabilities (0.8-1) with dark red.

393 We have distinguished two learning scenarios, one for validation (see sec-
394 tion 4.5) and the other for plotting the species distribution maps described
395 above. The models used to plot the maps were learnt from a random subsam-
396 ple containing an 80% of the original dataset, aiming at avoiding overfitting.

397 4.5. Model validation

398 The models were validated by means of the k -fold cross validation tech-
399 nique (Stone, 1974), obtaining a confusion matrix in each fold. This technique
400 randomly splits the complete dataset into k subsets, using $k-1$ of them for
401 learning the model and one for validation. The method is repeated k times
402 and the confusion matrix of the model is computed in each step. In order
403 to compute a confusion matrix, the output of the model, given in terms of
404 “probability of presence of the species”, P_S , was transformed into absence
405 and presence records, coded as “1” if $P_S \geq 0.5$ or “0” if $P_S < 0.5$ respectively.
406 In this case study, a k value of 10 was applied and therefore 10 confusion
407 matrices of each model were obtained.

408 A number of statistics were computed from the confusion matrices.
409 Their definition is given in Table 2. Of special interest is the area under
410 the ROC curve (AUC), which measures the predictive power of the model

411 considering both the true positive (Recall) and the true negative (Specificity)
412 rates. When AUC is defined by only 1 run, it is known as *balanced accuracy*
413 (Sokolova et al., 2006).

414 The models were validated in terms of Accuracy, Recall, f – score and
415 AUC. Precision and Specificity are also defined in Table 2 because they are
416 used in the definition of the four selected statistics. All of them measure
417 some aspect of the performance of the model and range from 0 to 1, with
418 models scoring close to 1 being better. In the case of AUC, a score of 0.5
419 represents a predictive power no better than predicting the class at random.
420 Wilcoxon’s signed rank Test was performed to detect differences between the
421 performance statistics of TAN and Zi-TAN models.

422 4.6. Results and discussion of the case study

423 4.6.1. Salamander

424 Figure 6 represents the qualitative component of the TAN (a) and Zi-TAN
425 (b) models for salamander. Figure 7 shows the box plots corresponding to
426 the 10 measures of the 4 performance statistics for TAN and Zi-TAN, along
427 with the p – values obtained in the pairwise comparison.

428 The Wilcoxon Test showed no significant differences ($p > 0.05$) between
429 both models regarding Accuracy. However, this measure is not recommend-
430 able when the distribution of the classes is unequal, as in this case where there
431 are 300 presences and 587 absences. On the other hand, the test showed sig-
432 nificant differences ($p < 0.05$) between TAN and Zi-TAN for the remaining
433 statistics: Recall, f-score and AUC. According to the statistical test applied,
434 Zi-TAN performs better than TAN with reference to the true presence rate
435 (Recall), i.e. the proportion of observed presences correctly classified was
436 higher in the Zi-TAN model. Regarding f-score and AUC measures, Zi-TAN
437 model scores higher than TAN, suggesting that the former classifies better
438 than the latter in the case of salamander.

439 Figure 8 shows the probability of presence of the salamander according
440 to the TAN (a) and Zi-TAN (b) models, along with its observed distribution
441 area. Both models recognize the general distribution pattern of the species,
442 identifying 3 main populations: 1 sub-population located in the north, along
443 the Sierra Morena mountain range; another sub-population located in the
444 northeast, in the Prebaetic System; and another sub-population located in
445 the south-southwest, in the Penibaetic System. Note that Zi-TAN misclassi-
446 fied more observed absences than TAN, i.e. the type I error is lower in the
447 TAN model. From the ecological point of view, the type I error or, in this

448 case, false presence may be understood as the potential distribution area for
449 the species, i.e., where salamander has never been seen but the environmen-
450 tal conditions meet its necessities. According to Martin et al. (2005) and
451 Lecomte et al. (2013), observed true absences may occur if the species does
452 not saturate its entire suitable area. The potential distribution area, based
453 on the selected variables, predicted by the Zi-TAN model nearly corresponds
454 to the observed distribution of the species.

455 4.6.2. Eagle

456 Figure 9 represents the qualitative component of the TAN (a) and Zi-
457 TAN (b) models for eagle. Figure 10 shows the box plots corresponding to
458 the 10 measures of the 4 performance statistics for TAN and Zi-TAN, along
459 with the p – values obtained in the pairwise comparison.

460 The Wilcoxon Test showed no significant differences ($p > 0.05$) between
461 both models regarding Accuracy. As mentioned above, Accuracy may not
462 be appropriate, especially for imbalanced datasets (Chawla, 2005). On the
463 other hand, the test showed significant differences ($p < 0.05$) between TAN
464 and Zi-TAN for Recall and AUC. The f-score statistic could not be calculated
465 for TAN since every fold of the cross validation yielded a Recall of 0 and an
466 undefined Precision (0 divided by 0) due to the fact that the model only
467 predicted absences in this case. According to the statistical test applied,
468 Zi-TAN performs better than TAN with reference to the true presence rate
469 (Recall), i.e. the proportion of observed presences correctly classified was
470 higher in the Zi-TAN model. Regarding AUC, Zi-TAN model scores higher
471 than TAN, suggesting that the former classifies better than the latter in the
472 case of eagle.

473 Figure 11 shows the potential distribution area of eagle, based on the
474 selected variables, given by TAN (a) and Zi-TAN (b) along with its observed
475 distribution area. Examining these maps, it is noticeable that the TAN clas-
476 sifier is more conservative than Zi-TAN, since the former relies more on the
477 probability of the dominant class, absences, and barely classifies observations
478 as presences. The TAN model obtained just 5 cells with probability of pres-
479 ence higher than zero. The poor model’s performance may be due to the
480 combination of a great number of zero-inflated explanatory variables (10 out
481 of 15) and the imbalanced class variable (55 observed presences out of 887
482 observations). In contrast, the Zi-TAN model fairly detected the 3 main pop-
483 ulations of eagle in Andalusia: Doñana, Eastern Sierra Morena and Central
484 Sierra Morena. The model also marked, with low probability, the Campo

485 de Gibraltar county, in the south, as a potential distribution area of eagle,
486 which was occupied by the eagle in the past (Gonzalez et al., 1989; González
487 et al., 2008).

488 **5. Conclusions**

489 We have developed Zi-TAN, a new model for dealing with *zero-inflated*
490 feature variables, using hybrid Bayesian networks. Our experimental re-
491 sults showed strong evidence that our proposed methodology for modeling
492 explanatory variables with zero excess improves the performance of the clas-
493 sifier. In the case of salamander, an abundant species in the study area,
494 the TAN model recognized the general pattern of the species and had a fair
495 performance whereas the distribution area predicted by the Zi-TAN model
496 corresponds almost exactly to the observed distribution area. In the case of
497 eagle, a scarce species in the study area, the TAN model had a poor perfor-
498 mance while Zi-TAN substantially improved the distribution area predicted
499 by the former. The technique explained in this paper can be applied to
500 species distribution models where the explanatory variables have an exces-
501 sive number of zeros. Further research needs to be done in order to argue its
502 application to other disciplines.

503 **Acknowledgements**

504 This work has been supported by the Spanish Ministry of Economy and
505 Competitiveness, through project TIN2013-46638-C3-1-P , by Junta de An-
506 dalucía through project P11-TIC-7821 and by ERDF-FEDER funds. A.D.
507 Maldonado is being supported by the Spanish Ministry of Education, Culture
508 and Sport through an FPU research grant, FPU2013/00547.

509 **References**

- 510 Aguilera, P. A., Fernández, A., Fernández, R., Rumí, R., Salmerón, A., 2011.
511 Bayesian networks in environmental modelling. *Environmental Modelling*
512 & *Software* 26, 1376–1388.
- 513 Aguilera, P. A., Fernández, A., Reche, F., Rumí, R., 2010. Hybrid Bayesian
514 network classifiers: Application to species distribution models. *Environ-*
515 *mental Modelling & Software* 25 (12), 1630–1639.

- 516 Aguilera, P. A., Fernández, A., Roperó, R. F., Molina, L., 2013. Ground-
517 water quality assessment using data clustering based on hybrid Bayesian
518 networks. *Stochastic Environmental Research & Risk Assessment* 27 (2),
519 435–447.
- 520 Ancelet, S., Etienne, M.-P., Benoît, H., Parent, E., 2010. Modelling spa-
521 tial zero-inflated continuous data with an exponentially compound Poisson
522 process. *Environmental and Ecological Statistics* 17, 347–376.
- 523 Böhning, D., Dietz, E., Schlattmann, P., 1999. The zero-inflated Poisson
524 model and the decayed, missing and filled teeth index in dental epidemi-
525 ology. *Journal of the Royal Statistical Society A* 162(2), 195–209.
- 526 Calama, R., Mutke, S., Tom’è, J., Gordo, J., Monterio, G., Tomé, M., 2011.
527 Modelling spatial and temporal variability in a zero-inflated variable: The
528 case of stone pine (*Pinus pinea* L.) cone production. *Ecological Modelling*
529 222, 606–618.
- 530 Chawla, N., 2005. Data mining for imbalanced datasets: An overview. In:
531 Maimon, O., Rokach, L. (Eds.), *Data Mining and Knowledge Discovery*
532 *Handbook*. Springer US, pp. 853–867.
- 533 Cobb, B. R., Shenoy, P. P., Rumí, R., 2006. Approximating probability den-
534 sity functions with mixtures of truncated exponentials. *Statistics and Com-
535 puting* 16, 293–308.
- 536 Cragg, J. G., 1971. Some statistical models for limited dependent variables
537 with application to the demand for durable goods. *Econometrica* 39(5),
538 829–844.
- 539 Damgaard, C., 2008. Modelling pin-point plant cover data along an environ-
540 mental gradient. *Ecological Modelling* 214, 404–410.
- 541 Davis, J., 1986. *Statistical and Data Analysis in Geology*. J. Wiley.
- 542 Degani, G., Warburg, M., 1978. Population structure and seasonal activity
543 of the adult *Salamandra salamandra* (L.) (Amphibia, Urodela, Salaman-
544 dridae) in Israel. *Journal of Herpetology* 12, 437–444.

- 545 Dorevitch, S., Doi, M., Hsu, F.-C., Lin, K.-T., Roberts, J. D., Liu, L. C.,
546 Gladding, R., Vannoy, E., Li, H., Javor, M., Scheff, P. A., 2011. A compar-
547 ison of rapid and conventional measures of indicator bacteria as predictors
548 of waterborne protozoan pathogen presence and density. *Journal of Envi-*
549 *ronmental Monitoring* 13, 2427–2435.
- 550 Edmeades, S., Smale, M., 2006. A trait-based model of the potential demand
551 for a genetically engineered food crop in a developing economy. *Agricultural*
552 *Economics* 35, 351–361.
- 553 Elvira Consortium, 2002. Elvira: An Environment for Creating and Using
554 Probabilistic Graphical Models. In: *Proceedings of the First European*
555 *Workshop on Probabilistic Graphical Models*. pp. 222–230.
556 URL <http://leo.ugr.es/elvira>
- 557 Fernández, A., Morales, M., Salmerón, A., 2007. Tree augmented naïve Bayes
558 for regression using mixtures of truncated exponentials: Applications to
559 higher education management. *IDA'07. Lecture Notes in Computer Science*
560 4723, 59–69.
- 561 Friedman, N., Geiger, D., Goldszmidt, M., 1997. Bayesian network classifiers.
562 *Machine Learning* 29, 131–163.
- 563 Fytilis, N., Rizzo, D. M., 2013. Coupling self-organizing maps with a Naïve
564 Bayesian classifier: Stream classification studies using multiple assessment
565 data. *Water Resources Reseach* 49, 7747–7762.
- 566 Gonzalez, L. M., Hiraldo, F., Delibes, M., Calderon, J., 1989. Reduction
567 in the range of the Spanish Imperial Eagle (*Aquila adalberti* Brem, 1861)
568 since AD 1850. *Journal of Biogeography* 16, 305–315.
- 569 González, L. M., Oria, J., Sánchez, R., Margalida, A., Aranda, A., Prada,
570 L., Caldera, J., Molina, J. I., 2008. Status and habitat changes in the en-
571 dangered Spanish Imperial Eagle *Aquila adalberti* population during 1974-
572 2004: implications for its recovery. *Bird Conservation International* 18,
573 242–259.
- 574 Greene, W. H., 1994. Accounting for excess zeros and sample selection in
575 Poisson and Negative Binomial regression models. Tech. rep., Department
576 of Economics, Stern School of Business, New York University.

- 577 Hall, D. B., 2000. Zero-Inflated Poisson and Binomial regression with random
578 effects: A case study. *Biometrics* 56, 1030–1039.
- 579 Kamarianakis, Y., Feidas, H., Kokolatos, G., Chrysoulakis, N., Karatzias, V.,
580 2008. Evaluating remotely sensed rainfall estimates using nonlinear mixed
581 models and geographically weighted regression. *Environmental Modelling*
582 & *Software* 23, 1438–1447.
- 583 Lambert, D., 1992. Zero-Inflated Poisson regression, with an application to
584 defects in manufacturing. *Technometrics* 34, 1–14.
- 585 Langseth, H., Nielsen, T., Pérez-Bernabé, I., Salmerón, A., 2014. Learning
586 mixtures of truncated basis functions from data. *International Journal of*
587 *Approximate Reasoning* 55, 940–956.
- 588 Langseth, H., Nielsen, T. D., Rumí, R., Salmerón, A., 2009. Maximum like-
589 lihood learning of conditional MTE distributions. *ECSQARU’09. Lecture*
590 *Notes in Artificial Intelligence* 5590, 240–251.
- 591 Langseth, H., Nielsen, T. D., Rumí, R., Salmerón, A., 2012. Mixtures of
592 Truncated Basis Functions. *International Journal of Approximate Reason-*
593 *ing* 53 (2), 212–227.
- 594 Lecomte, J., Benoît, H., Etienne, M., Bel, L., Parent, E., 2013. Modeling the
595 habitat associations and spatial distribution of benthic macroinvertebrates:
596 A hierarchical Bayesian model for zero-inflated biomass data. *Ecological*
597 *Modelling* 265, 74–84.
- 598 Maldonado, A. D., Aguilera, P. A., Salmerón, A., 2015. Continuous Bayesian
599 networks for probabilistic environmental risk mapping. *Stochastic Environ-*
600 *mental Research & Risk Assessment* In press.
- 601 Markus, M., Hejazi, M. I., Bajcsy, P., Giustolisi, O., Savic, D. A., 2010. Pre-
602 diction of weekly nitrate-N fluctuations in a small agricultural watershed
603 in Illinois. *Journal of Hydroinformatics* 12.3, 251–261.
- 604 Martin, T. G., Wintle, B. A., Rhodes, J. R., Kuhnert, P. M., Field, S. A.,
605 Low-Choy, S. J., Tyre, A. J., Possingham, H. P., 2005. Zero tolerance
606 ecology: improving ecological inference by modelling the source of zero
607 observations. *Ecology Letters* 8, 1235–1246.

- 608 McDavid, A., Finak, G., Chattopadhyay, P. K., Dominguez, M., Ma, L. L.
609 S. S., Roederer, M., Gottardo, R., 2013. Data exploitation, quality control
610 and testing in single-cell qPCR-based gene expression experiments.
611 *Bioinformatics* 29(4), 461–467.
- 612 Moral, S., Rumí, R., Salmerón, A., 2001. Mixtures of Truncated Exponentials
613 in Hybrid Bayesian Networks. In: Benferhat, S., Besnard, P. (Eds.), *Symbolic and Quantitative Approaches to Reasoning with Uncertainty*. Vol.
614 2143 of *Lecture Notes in Artificial Intelligence*. Springer, pp. 156–167.
615
- 616 Mullahy, J., 1986. Specification and testing of some modified count data
617 models. *Journal of Econometrics* 33, 341–365.
- 618 Neil, M., Tailor, M., Marquez, D., 2007. Inference in hybrid Bayesian networks
619 using dynamic discretization. *Statistics and Computing* 17, 219–233.
- 620 Ngatchou-Wandji, J., Paris, C., 2011. On the zero-inflated count models with
621 application to modelling annual trends in incidences of some occupational
622 allergic diseases in France. *Journal of Data Science* 9, 639–659.
- 623 Nie, L., Wu, G., Brockman, F., Zhang, W., 2006. Integrated analysis of
624 transcriptomic and proteomic data of *Desulfovibrio vulgaris*: zero-inflated
625 Poisson regression models to predict abundance of undetected proteins.
626 *Bioinformatics* 22(13), 1641–1647.
- 627 Pearl, J., 1988. *Probabilistic reasoning in intelligent systems*. Morgan-
628 Kaufmann (San Mateo).
- 629 Potts, J. M., Elith, J., 2006. Comparing species abundance models. *Ecological*
630 *Modelling* 199, 153–163.
- 631 Ropero, R. F., Aguilera, P. A., Fernández, A., Rumí, R., 2014. Regression
632 using hybrid Bayesian networks: Modelling landscape-socioeconomy relationships.
633 *Environmental Modelling & Software* 54, 127–137.
- 634 Ropero, R. F., Aguilera, P. A., Rumí, R., 2015. Analysis of the socioecological
635 structure and dynamics of the territory using a hybrid Bayesian network
636 classifier. *Ecological Modelling* 311, 73–87.
- 637 Rumí, R., Salmerón, A., 2007. Approximate probability propagation with
638 mixtures of truncated exponentials. *International Journal of Approximate*
639 *Reasoning* 45, 191–210.

- 640 Shenoy, P. P., West, J. C., 2011. Inference in hybrid Bayesian networks using
641 mixtures of polynomials. *International Journal of Approximate Reasoning*
642 52 (5), 641–657.
- 643 Sokolova, M., Japkowicz, N., Szpakowicz, S., 2006. Beyond accuracy, f-score
644 and roc: A family of discriminant measures for performance evaluation.
645 In: Sattar, A., Kang, B.-h. (Eds.), *AI 2006: Advances in Artificial Intel-*
646 *ligence*. Vol. 4304 of *Lecture Notes in Computer Science*. Springer Berlin
647 Heidelberg, pp. 1015–1021.
- 648 Solé-Auró, A., Guillén, M., Crimmins, E. M., 2012. Health care usage among
649 immigrants and native-born elderly populations in eleven European coun-
650 tries: results from SHARE. *European Journal of Health Economics* 13,
651 741–754.
- 652 Stone, M., 1974. Cross-validatory choice and assessment of statistical predic-
653 tions. *Journal of the Royal Statistical Society. Series B (Methodological)*
654 36 (2), 111–147.
- 655 Varona, L., Sorensen, D., 2010. A genetic analysis of mortality in pigs. *Ge-*
656 *netics Society of America* 184, 277–284.
- 657 Wenger, S. J., Freeman, M. C., 2008. Estimating species occurrence, abun-
658 dance, and detection probability using zero-inflated distributions. *Ecology*
659 89(10), 2953–2959.

Table 1: Summary of variables. Note that not all these variables were included in the models but a variable selection process was carried out.

Variable	Description	Source
Salamander /Eagle	Presence/absence of the given species in each cell	Spanish inventory of terrestrial species ^a
T (°C)	Average of annual mean temperature for the 30 year period 1971-2000 in each cell	Annual Mean Temperature Dataset of Andalusia. TIFF raster format with 100 m spatial resolution ^b
Rainfall (mm)	Average of annual rainfall for the 30 year period 1971-2000 in each cell	Annual Precipitation Dataset of Andalusia. TIFF raster format with 100 m spatial resolution ^b
PET (mm)	Average of the annual potential evapotranspiration for the 30 year period 1971-2000 in each cell	Annual Mean Evapotranspiration Dataset of Andalusia. TIFF raster format with 100 m spatial resolution ^b
Humidity index	Average of annual humidity index for the 30 year period 1971-2000 in each cell	Annual Mean Humidity Index Dataset of Andalusia. Shapefile format ^b
Land uses (%)	Percentage of occupation of each land-use (#44) within each cell	Andalusian Land Use and Land Cover Map (1:25,000) ^b
Soil (%)	Percentage of occupation of each soil type (#63) within each cell	Andalusian Soil Map (1:400,000) ^b
Lithology (%)	Percentage of occupation of each lithological unit (#41) within each cell	Andalusian Lithological Map (1:400,000) ^b
Z (m a.s.l.)	Average elevation of each cell	Andalusian Digital Terrain Model. Grid width 200 m spatial resolution ^c
Slope (%)	Average slope of each cell	
Aspect (°)	Average aspect of each cell	

number of variables

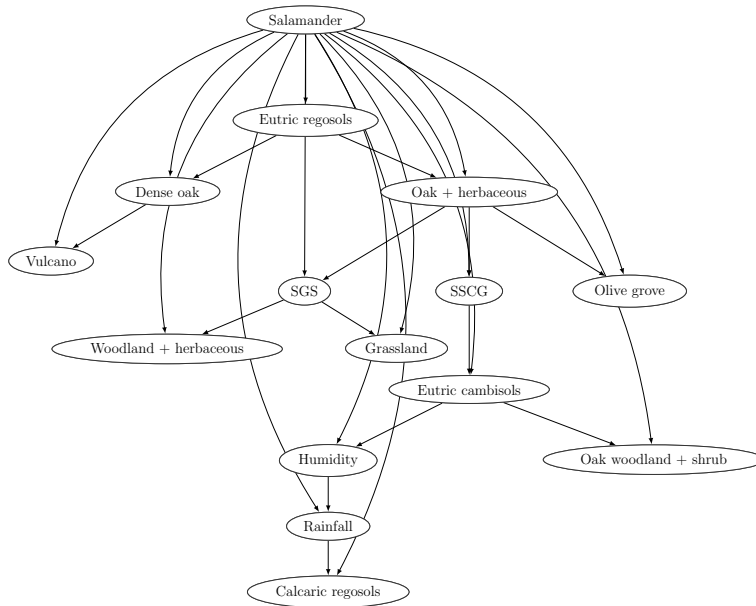
^aProvided by the Ministry of Agriculture, Food and Environment

^bProvided by the Andalusian Environmental Information Network

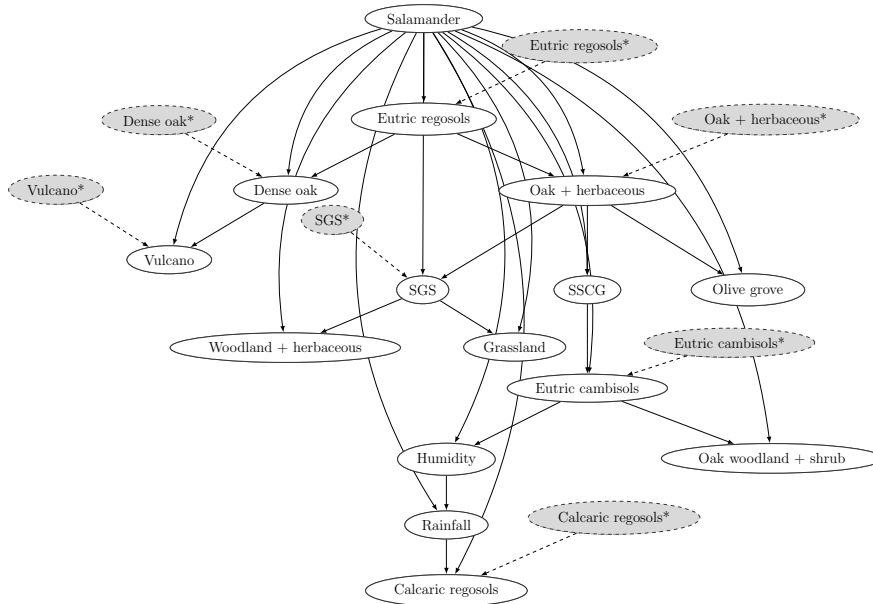
^cProvided by the Spanish National Geographic Institute

Table 2: Statistics used to validate the classification models. TP, TN, FP and FN are, respectively, the true positive, true negative, false positive and false negative rates. Parameter β is the relative importance of Precision vs Recall, and was set to 1 in the case study.

Statistic	Definition	Statistic	Definition
Accuracy	$\frac{TP+TN}{TP+FN+FP+TN}$	f - score	$\frac{(1+\beta) \times \text{Recall} \times \text{Precision}}{\beta^2 \times \text{Recall} + \text{Precision}}$
Recall	$\frac{TP}{TP+FN}$	Specificity	$\frac{TN}{TP+FN}$
Precision	$\frac{TP}{TP+FP}$	AUC	$\frac{1}{2}(\text{Recall} + \text{Specificity})$



(a) TAN structure



(b) Zi-TAN structure

Figure 6: Structure of TAN (a) and Zi-TAN (b) models for predicting presence of salamander given the explanatory variables. SGS: slate-greywacke-sandstone; SSCG: sand-silt-clay-gravel. Shaded nodes with dashed lines represent the artificial binary variables used to model the zero-inflated explanatory variables. Note that the dependence relationships existing between the feature variables just allow the models to perform better.

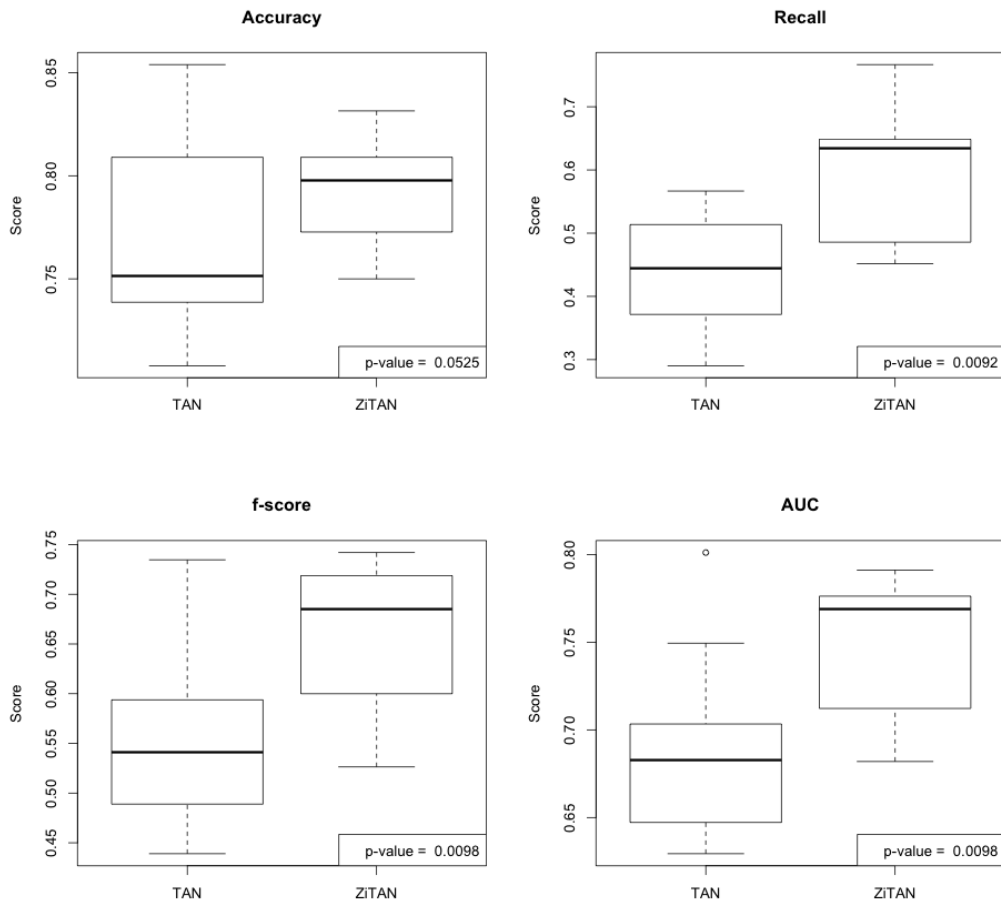


Figure 7: Box plots comparing TAN and Zi-TAN in terms of their performance statistics in the *Salamandra salamandra* case. The p - values obtained in the Wilcoxon Test are shown.

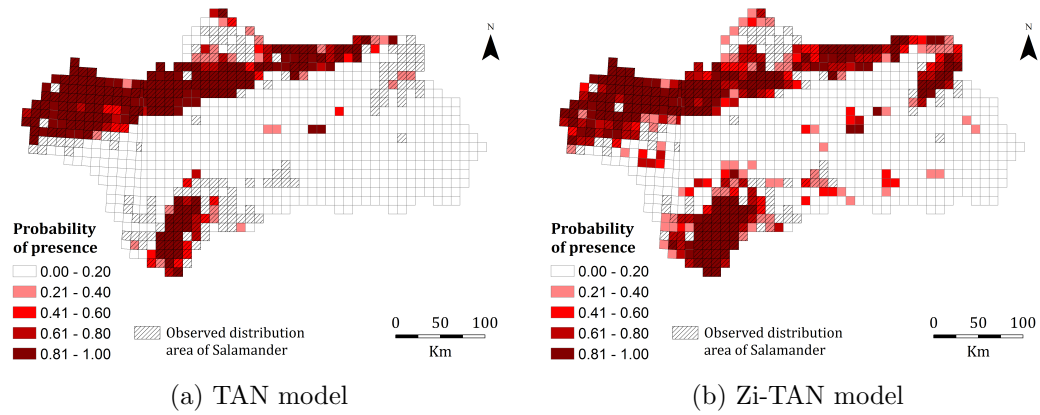
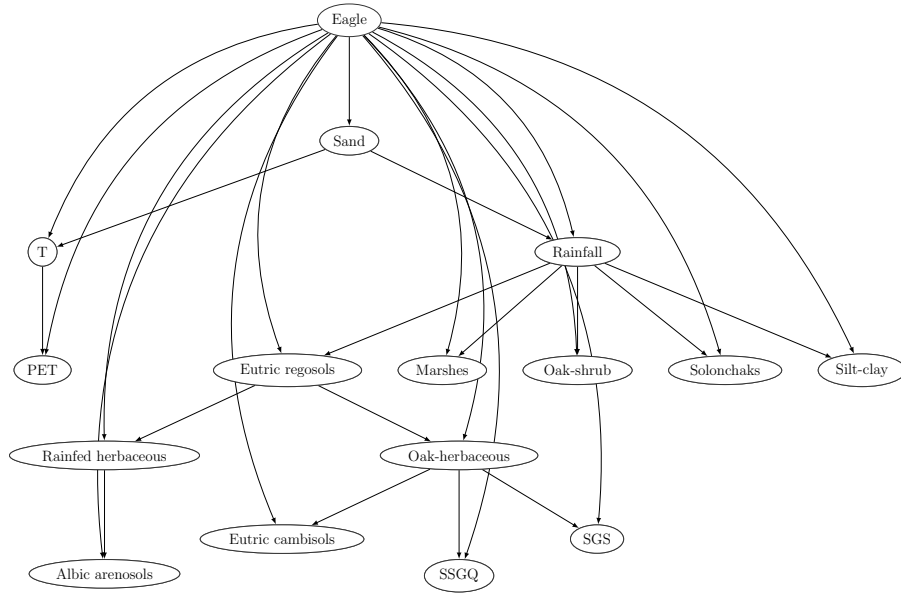
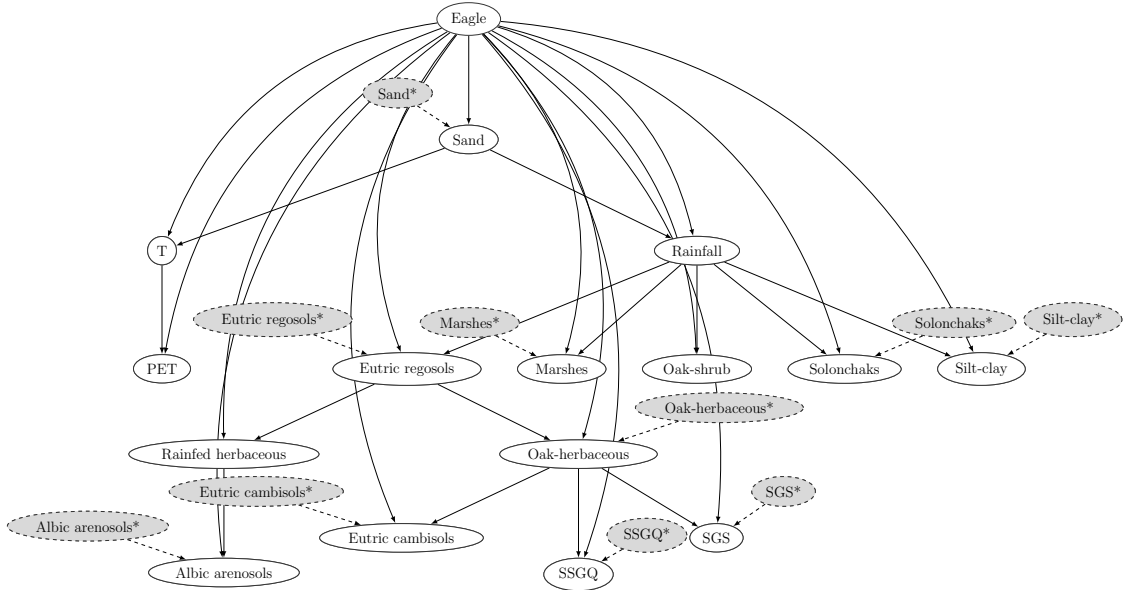


Figure 8: Potential distribution area of salamander, expressed as probability of presence, predicted by TAN (a) and Zi-TAN (b) models. Cells filled in with straight lines represent the observed distribution area of the species.



(a) TAN structure



(b) Zi-TAN structure

Figure 9: Structure of TAN (a) and Zi-TAN (b) models for predicting presence of eagle given the explanatory variables. SSGQ: slate-shale-greywacke-quartzite; SGS: slate-greywacke-sandstone. Shaded nodes with dashed lines represent the artificial binary variables used to model the zero-inflated explanatory variables. Note that the dependence relationships existing between the feature variables just allow the models to perform better.

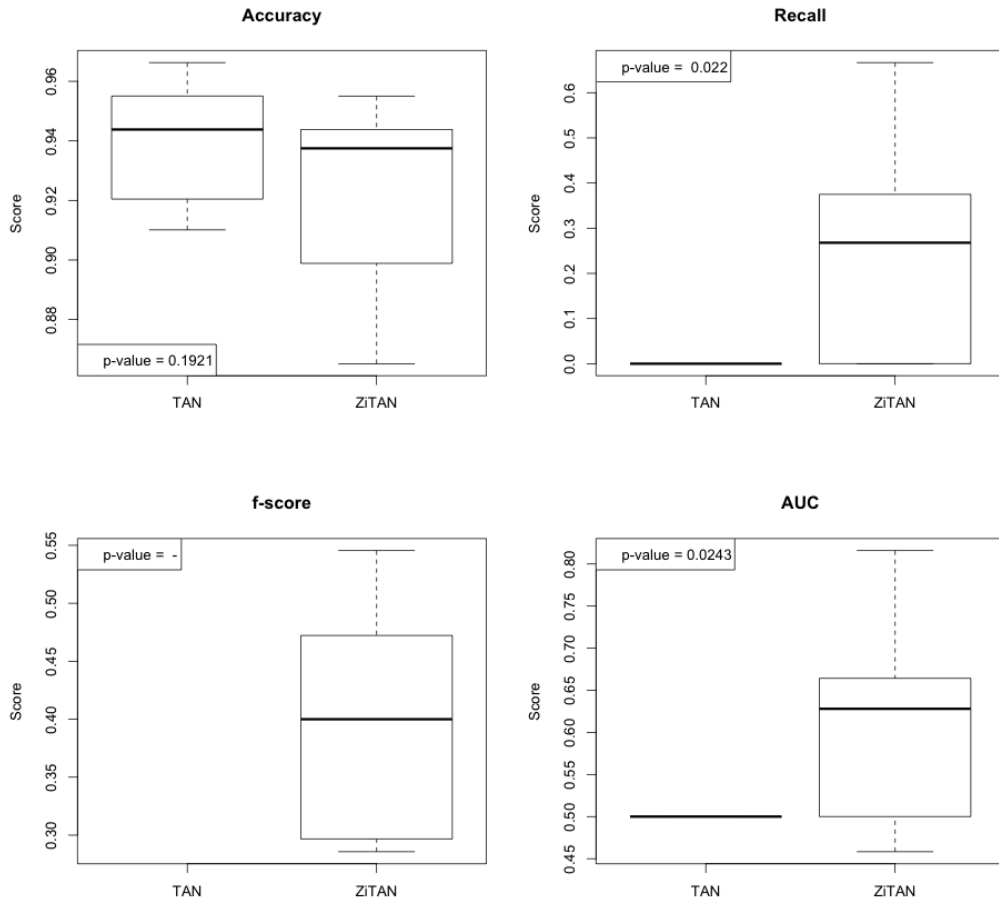


Figure 10: Box plots comparing TAN and Zi-TAN in terms of their performance statistics in the *Aquila adalberti* case. The p -values obtained in the Wilcoxon Test are shown.

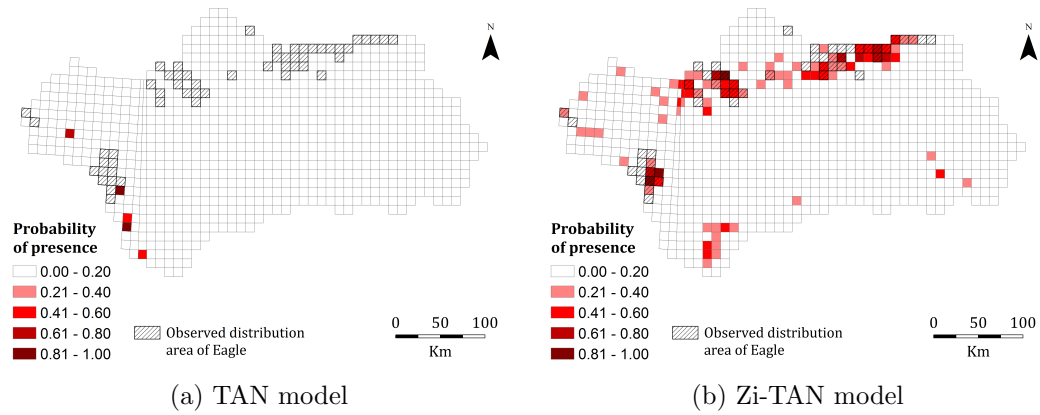


Figure 11: Potential distribution area of eagle, expressed as probability of presence, predicted by TAN (a) and Zi-TAN (b) models. Cells filled in with straight lines represent the observed distribution area of the species.