

# Discretizing Environmental Data for Learning Bayesian-network Classifiers

R.F. Ropero<sup>a</sup>, S. Renooij<sup>b</sup>, L.C. van der Gaag<sup>b</sup>

<sup>a</sup>*Informatics and Environment Laboratory, Dept. of Biology and Geology, University of Almería, Carretera de Sacramento s/n, La Cañada de San Urbano, Almería, Spain*

<sup>b</sup>*Department of Information and Computing Sciences, Utrecht University, Princetonplein 5, De Uithof, Utrecht, The Netherlands*

---

## Abstract

For predicting the presence of different bird species in Andalusia from land-use data, we compare the performances of Bayesian-network classifiers and logistic-regression models. In our study, both balanced and unbalanced data sets are used, and models are learned from both the original continuous data and from the data after discretization. For the latter purpose, four different discretization methods, called *Equal Frequency*, *Equal Width*, *Chi-Merge* and *MDLP*, are compared. The experimental results from our species data sets suggest that the simple Naive Bayesian classifiers are preferable to logistic-regression models and that the relatively unknown *Chi-Merge* method is the preferred method for discretizing these environmental data.

*Keywords:* Species distribution models, Bayesian-network classifiers, Logistic-regression models, Discretization methods

---

## 1. Introduction

Bayesian networks (BNs for short) are powerful probabilistic models that have demonstrated their usefulness in a wide range of application fields among which is the environmental-science field (Baur and Bozdog, 2015; Jensen and Nielsen, 2007). In environmental science, Bayesian networks are used for knowledge discovery, where the focus is on establishing the relationships among the variables at hand and their evolution under various scenarios

---

*Email addresses:* `rosa.ropero@ual.es` \*\*Corresponding author (R.F. Ropero), `S.Renooij@uu.nl` (S. Renooij), `L.C.vanderGaag@uu.nl` (L.C. van der Gaag)

8 (Dyer et al., 2014). Bayesian networks are further used for classification pur-  
9 poses (Maldonado et al., 2015; Park and Stenstrom, 2008), where the aim is  
10 to accurately predict the value of a specific target variable, called the class  
11 variable.

12 Initially, Bayesian networks were designed to handle data pertaining to  
13 discrete variables only. Real-world data are often of a continuous or hybrid  
14 nature however, and new algorithms for learning and inference in Bayesian  
15 networks with both continuous and discrete variables are emerging (Langseth  
16 et al., 2012; Moral et al., 2001). Despite the increasing availability of such  
17 algorithms, most Bayesian-network packages to date require variables to be  
18 discrete. Upon practical application, therefore, any continuous variables need  
19 to be discretized.

20 Discretization is widely applied in knowledge-discovery and machine-  
21 learning applications, with the aim of *i*) reducing and simplifying the avail-  
22 able data, *ii*) rendering model learning more efficient, and *iii*) obtaining more  
23 compact and more readily interpretable results (Liu et al., 2002). Over the  
24 years, several different discretization methods have been proposed, only a few  
25 of which are widely used while others are largely unnoticed (García et al.,  
26 2013; Yang et al., 2010; Liu et al., 2002). Since data discretization gener-  
27 ally results in information loss (Li, 2007; Uusitalo, 2007), the discretization  
28 method employed will affect the predictive quality of any model learned from  
29 the data. Where several papers address the question of which discretization  
30 method is most suited for data mining in general (García et al., 2013; Liu  
31 et al., 2002) or for Bayesian-network learning in particular (Lima et al., 2014;  
32 Zhou et al., 2014), the best choice of method tends to depend on the nature  
33 and characteristics of the data at hand.

34 In environmental science, Bayesian networks are typically used in a decision-  
35 making process in which expert knowledge plays an important role (Voinov  
36 and Bousquet, 2010). In this context, the use of discrete data provides more  
37 easily interpretable results and facilitates the communication between mod-  
38 elers and environmental experts (García et al., 2013; Liu et al., 2002). Ac-  
39 cording to a recent review (Aguilera et al., 2011), in fact, more than 80%  
40 of the papers addressing Bayesian networks in environmental science involve  
41 discretized data, where the discretization is done using the so-called *Equal*  
42 *Frequency* method or is based on expert knowledge. While more tailored  
43 discretization methods have been designed for specific types of model, such  
44 as hydrological models (Pradhanang and Briggs, 2014), models of air qual-  
45 ity (Davison and Ramesh, 1996), and models of spatial distributions of the

46 data (Liu et al., 2015), discretization methods specifically designed for en-  
47 vironmental modeling through Bayesian networks do not abound. To bring  
48 the discretization methods in use with Bayesian networks in general to the  
49 attention of environmental modelers, further efforts as well as more tailored  
50 insights are called for (Nash et al., 2013).

51 During the last decades, species distribution modeling has evolved in  
52 the field of environmental science, following the development of Geographic  
53 Information Systems (GIS) and spatial statistics techniques (Segurado and  
54 Araújo, 2004). In general, the objective of species distribution modeling is  
55 to link species data with environmental variables and to obtain maps show-  
56 ing the spatial distribution of the species under study (Elith et al., 2006).  
57 Some of the most commonly used models for this purpose are classification  
58 trees (Fukuda et al., 2013), regression models (Li and Wang, 2013), neural  
59 networks (Dedecker et al., 2004), and more tailored models like BIOCLIM  
60 (Busby, 1986) and FLORAMAP (Jones and Gladkov, 1999). In contrast,  
61 Bayesian networks are scarcely being applied in species distribution model-  
62 ing, although some examples are found, addressing classification with dis-  
63 cretized data (Newton et al., 2007) and using a model structure based on  
64 expert knowledge (Pollino et al., 2007).

65 In this paper we compare various classification models for predicting the  
66 presence of different bird species in Andalusia from land-use data. More  
67 specifically, we study the performance of two types of Bayesian-network clas-  
68 sifier: the Naive Bayesian (NB) classifier and the Tree Augmented Naive  
69 Bayesian (TAN) classifier. These classifiers are learned from both the original  
70 continuous data and from discretized data. For discretization, four methods  
71 are compared: *Equal Frequency* (EF), *Equal Width* (EW), *Chi-Merge* (ChiM)  
72 and a method based on the *Minimum Description Length Principle* (MDLP);  
73 these methods are the most commonly used discretization methods (García  
74 et al., 2013; Liu et al., 2002). We further compare the performances of these  
75 classifiers when learned from well-balanced data sets and from less balanced  
76 data.

77 The performance of a classification model depends to a large extent on  
78 the decision rule that is used to decide upon the class to which a case is  
79 assigned. In practice often maximum-probability classification is used, in  
80 which a case is assigned to the most likely class (Ropero et al., 2015; Aguil-  
81 era et al., 2013). In essence, however, any probability can be chosen for a  
82 decision threshold: a species then is classified as *present* if the predicted prob-  
83 ability of it being present exceeds this threshold, and as *absent* otherwise.

84 For less balanced data sets, in which the prior distribution over the class  
85 variable is quite skewed, maximum-probability classification may lead to un-  
86 desirable classification behaviour (van der Gaag et al., 2009a). In this paper  
87 we therefore study the performance of the various classifiers with maximum-  
88 probability classification and with threshold-probability classification using  
89 a decision threshold based on the prior species distribution (van der Gaag  
90 et al., 2009b).

91 Since in species distribution modeling the use of logistic-regression models  
92 is quite common, from the various data sets also logistic-regression models  
93 are constructed and compared with the learned Bayesian-network classifiers  
94 in terms of their performance.

## 95 **2. Materials and Methods**

96 In this section we review the data sets available for our study and describe  
97 the various methods used for discretizing these data and for learning and  
98 validating classification models.

### 99 *2.1. Study area and data collection*

100 Andalusia, located in the South of Spain (Fig. 1), constitutes the na-  
101 tion's second largest autonomous region, with a surface area of 87,600 km<sup>2</sup>  
102 representing 17.3% of the national territory<sup>1</sup>. Lying on the frontier between  
103 Europe and Africa, Andalusia inherits landscape and biodiversity specifics  
104 from both continents. Its terrain covers a wide range of altitudes, from the  
105 Baetic Depression to the mountainous ranges of the Sierra Morena and the  
106 Baetic System, with the highest peaks lying over 3000 meters above sea level  
107 (m.a.s.l.) The landscape is quite heterogeneous, with huge differences from  
108 the densely populated and irrigated cropland areas of the river basin and  
109 coastlands, to the sparsely populated forested areas of the uplands. Its cli-  
110 mate is similarly heterogeneous, with stark differences between inland and  
111 coastal areas. The climate in the south-eastern coastal part is semiarid, with  
112 less than 200 mm of annual rainfall in several areas, while the middle and  
113 northern parts have a continental climate, with more than 4000 mm of rainfall  
114 per year. These natural conditions make Andalusia a heterogeneous region

---

<sup>1</sup>Data from the Spanish Statistical Institute.

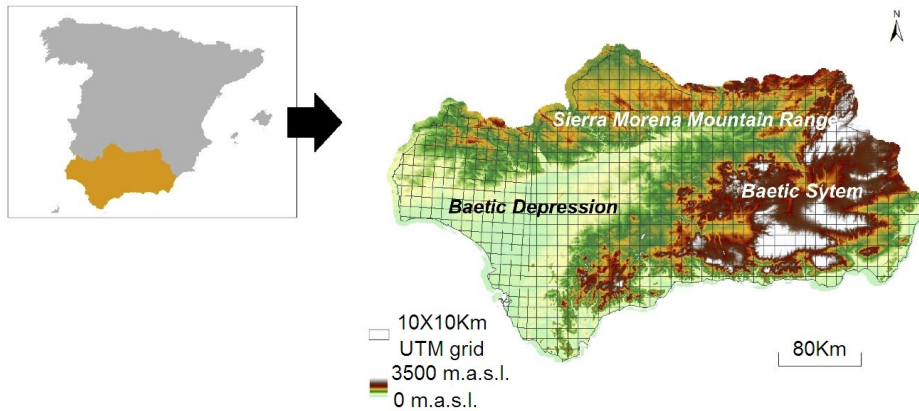


Figure 1: Andalusia, located in the South of Spain (*left*), its relief and the UTM  $10 \times 10$  km grid used for the data collection (*right*); the smaller cells in the western area result from the grid having been corrected to fit two geographical HUSOS.

115 both in terms of territorial structure and in climatic and ecological condi-  
 116 tions. Provoking ecological niches with large biodiversity rates, Andalusia is  
 117 considered a global biodiversity hotspot (Myers et al., 2000).

118 The Spanish Inventory of Terrestrial Species<sup>2</sup> by the Spanish National  
 119 Government was used to select information about the prevalence of three  
 120 bird species – *Turdus viscivorus*, *Cecropis daurica* and *Accipiter nisus* – for  
 121 the UTM (*Universal Transverse Mercator*)  $10 \times 10$  km grid of Andalusia  
 122 (Fig. 1); the three species were selected for their different prevalence rates.  
 123 Information about land use for the same grid was collected from the An-  
 124 dalusian Environmental Information Network<sup>3</sup> from the Andalusian Regional  
 125 Government. ArcGIS 9.3 was used for selecting the data and merging them  
 126 into the grid. As a consequence of the high heterogeneity of the region, a  
 127 single cell of the grid of Andalusia can show several small patches of different  
 128 types of land use, as illustrated in Fig. 2(a). A more detailed example, show-  
 129 ing the distribution of *Olive cropland* coverage over the grid cells, is provided  
 130 in Fig. 2(b). The figure shows that, for this land-use variable, the majority

<sup>2</sup><http://www.magrama.gob.es/es/biodiversidad/temas/inventarios-nacionales/inventario-especies-terrestres/>

<sup>3</sup><http://www.juntadeandalucia.es/medioambiente/site/rediam>

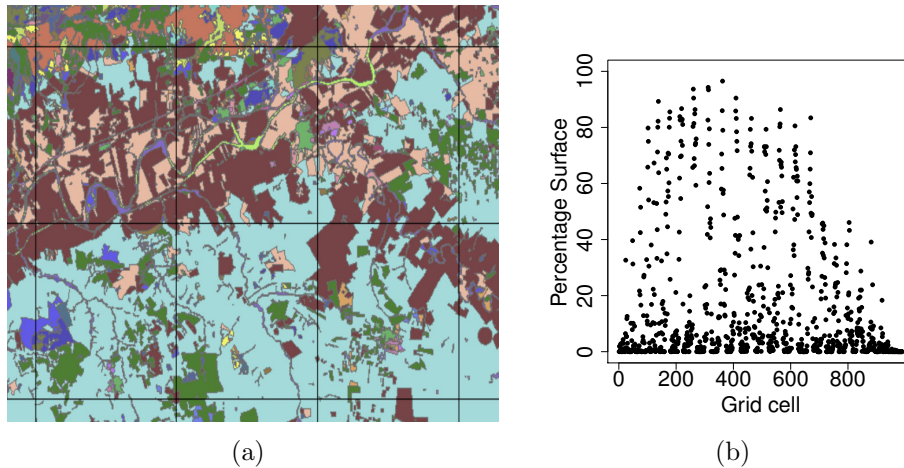


Figure 2: Enlarged part of the  $10 \times 10$  km grid showing different types of land use (a), and the distribution of *Olive cropland* coverage over all grid cells.

131 of recorded percentages are within the range of 0% to 10% of the surface,  
 132 while the remaining data values are scattered over the 10% to 100% interval.  
 133 Similarly skewed distributions are found for all variables involved.

134 The data used for our study is composed of three data sets, one for each  
 135 bird species of interest. Each data set includes a single discrete class variable  
 136 that represents whether the bird species at hand is *present* or *absent* in a  
 137 specific grid cell. The remaining variables, listed in Table 1, are continuous  
 138 feature variables which represent the percentage (between 0% and 100%) of  
 139 a grid cell’s surface with a particular type of land use. The actual features  
 140 were extracted from regional reports about the biology of each species, by  
 141 selecting those pertaining to the species’ habitat. Each data set contains 989  
 142 records, one per grid cell, and does not have any missing values, that is, for  
 143 each grid cell, both the actual features and the associated class are recorded.

## 144 2.2. Classification models

145 Two types of Bayesian-network classifier are studied, each with discretized  
 146 variables and with the original continuous variables respectively, and their  
 147 performances are compared with those of a logistic-regression model.

### 148 2.2.1. Bayesian-network classifiers

149 A Bayesian network is a concise model of a joint probability distribution  
 150 over a set of random variables (Jensen and Nielsen, 2007). It combines a

Table 1: Feature variables and prevalence ( $p$ ) per species.

	<i>Turdus viscivorus</i> $p = 0.47$	<i>Cecropis daurica</i> $p = 0.84$	<i>Accipiter nisus</i> $p = 0.27$
1	Bare soil	Agricultural areas	Bare soil
2	Dams	Bare soil	Bare soil of scrub
3	Dense forest of conifers	Cliff	Dense forest of conifers
4	Dense forest of oaks	Dehesas	Dense forest of oaks
5	Dense scrubland	Dense forest	Dense grasslands
6	Dense scrubland of conifers	Dense scrubland	Dense scrubland
7	Dense scrubland of oaks	Dense scrubland of trees	Dense scrubland of oaks
8	Open scrubland	Open scrubland	Open grasslands
9	Grasslands of oaks	Open scrubland of trees	Open scrubland
10	Herbaceous crops	Grasslands	Open scrubland of oaks
11	Heterogeneous crops	Grasslands of trees	Grasslands of oaks
12	Irrigation pond	Herbaceous crops	Herbaceous crops
13	Olive crops	Heterogeneous crops	Heterogeneous crops
14	Other dense forests	Man-made water surfaces	Irrigated pool
15	Woody crops	River bed	Olive crops
16		Urban areas	Other disperse scrubland of trees
17		Woody crops	Other dense forest
18			Other dense scrubland of trees
19			River bed
20			Woody crops

151 directed acyclic graph, which describes the (in)dependencies between the  
 152 variables, with local probability distributions per variable. From a Bayesian  
 153 network, any probability of interest over its variables can be computed.

154 When used for classification purposes, a Bayesian network includes a  
 155 designated class variable  $C$ . Of interest then is the posterior probability dis-  
 156 tribution  $\Pr(C \mid \mathbf{f})$  over  $C$  given case observations  $\mathbf{f}$  for the feature variables  
 157 involved. To decide upon the class to which the observations  $\mathbf{f}$  are to be  
 158 assigned, two approaches are in use:

- 159 • *maximum-probability classification* (also known as “the winner takes  
 160 all”), in which the case observations  $\mathbf{f}$  are assigned to the most probable  
 161 class given  $\mathbf{f}$ ;

- 162 • *probability-threshold classification*, in which the observations  $\mathbf{f}$  are as-  
163 signed to the class  $c$  if  $t_1 > \Pr(c | \mathbf{f}) \geq t_2$  for some suitable choice of  
164 decision thresholds  $t_1$  and  $t_2$ .

165 For a binary class variable with the classes  $c_1$  and  $c_2$ , probability-threshold  
166 classification with a decision threshold  $t$  assigns case observations  $\mathbf{f}$  to  $c_1$  if

$$\Pr(c_1 | \mathbf{f}) \geq t$$

167 and to  $c_2$  otherwise; taking  $t = 0.5$  would then result in the same class  
168 assignment as maximum-probability classification. The overall performance  
169 of a probabilistic classifier is optimized by choosing a decision threshold based  
170 on the prior distribution over the class variable (Lachiche and Flach, 2003).

171 For classification purposes, tailored Bayesian networks with highly con-  
172 strained graphical structures are in use, among which are the Naive Bayesian  
173 (NB) classifier and the Tree Augmented Naive Bayesian (TAN) classifier  
174 (Friedman et al., 1997). The Naive Bayesian classifier is the most constrained  
175 of all Bayesian-network classifiers: its graph consists of a designated node for  
176 the class variable and nodes modeling the feature variables with just this  
177 class variable for their parent. This type of classifier derives its name from  
178 the fact that its graphical structure captures the naive assumption that all  
179 feature variables are mutually independent given the class variable. Although  
180 this assumption does not generally hold in practice, NB classifiers tend to  
181 show quite competitive performance. TAN classifiers allow for explicitly rep-  
182 resenting dependencies among the feature variables by a tree structure, and  
183 in essence may thereby outperform NB classifiers.

184 Learning a Naive Bayesian classifier from a data set amounts to estimat-  
185 ing probabilities from the available data so as to quantify the relationships  
186 between the class variable and each of the feature variables. Learning a TAN  
187 classifier in addition involves learning the graphical structure from the data.  
188 For this purpose, first a directed tree over the feature variables is learned by  
189 building upon the conditional mutual information between pairs of feature  
190 variables given the class variable (Chow and Liu, 1968); subsequently, the  
191 class variable is added and all modeled relationships are quantified.

### 192 2.2.2. *Hybrid Bayesian networks*

193 Bayesian networks were initially defined for discrete variables only. Even  
194 to date, Bayesian-network software packages tend to assume all variables to  
195 be discrete. As a consequence, upon developing a real-world application, all



196 continuous domain variables have to be discretized by dividing their value  
197 ranges into a sequence of adjacent intervals. A probability distribution over  
198 the discretized variable then assigns to each such interval a single probability  
199 which can be viewed as approximating the continuous distribution over the  
200 interval by a constant function. In general, the use of more intervals upon  
201 discretization tends to result in a better approximation, albeit at the expense  
202 of a more complex model.

203 More recently, approaches have been developed that allow Bayesian net-  
204 works to include both continuous and discrete variables (Langseth et al.,  
205 2012; Shenoy and West, 2011; Lauritzen and Jensen, 2001; Moral et al., 2001).  
206 In this paper we study Bayesian-network classifiers that employ *Mixtures of*  
207 *Truncated Exponentials* (MTEs) for their local probability distributions. Like  
208 discretization methods, MTE approaches divide the value range of a contin-  
209 uous variable into intervals. The continuous distribution per interval is then  
210 approximated by an exponential function rather than by a constant function  
211 (Rumí, 2003). Similar to discretization, the use of more intervals tends to  
212 result in a better approximation, but will also yield a more complex model.  
213 By including more terms in the MTE per interval, the approximation also  
214 tends to improve, yet again at the cost of a more complex model (Rumí and  
215 Salmerón, 2007; Morales et al., 2006; Rumí et al., 2006).

### 216 2.2.3. Logistic regression

217 Logistic regression is a type of regression in which a binary response vari-  
218 able (the binary class variable, in terms of our classification context) is related  
219 to multiple explanatory variables (the feature variables) which may be dis-  
220 crete or continuous (Scott, 2010). Upon classification of case observations  $\mathbf{f}$   
221 for the explanatory variables, the response with highest posterior odds given  
222  $\mathbf{f}$  is determined and assigned to the case. Logistic-regression classification  
223 thereby in essence is similar to taking a maximum-probability approach to  
224 classification. In fact, logistic-regression classification is known to be equiv-  
225 alent to Naive-Bayesian classification under mild conditions (Roos et al.,  
226 2005).

### 227 2.3. Data discretization

228 In our study, four discretization methods are compared: *Equal Frequency*,  
229 *Equal Width*, *Chi-Merge* and *Minimum Description Length Principle* dis-  
230 cretization. These four methods are the most commonly studied discretiza-  
231 tion methods in the literature (García et al., 2013; Liu et al., 2002). In

232 environmental modeling with Bayesian networks, the *Equal Frequency* and  
233 *Equal Width* methods prevail (Aguilera et al., 2011). The *Chi-Merge* method  
234 has, to the best of our knowledge, never been used in such applications, while  
235 use of the *Minimum Description Length Principle* has been reported in just  
236 a single environmental-modeling study (Fernandes et al., 2013).

### 237 2.3.1. Discretization methods

238 Application of a discretization method to a data set starts by sorting the  
239 available data points in increasing order of their value for the continuous  
240 variable to be discretized. The data points are then distributed over  $k > 1$   
241 bins, each of which is associated with an interval from the variable’s overall  
242 value range. Discretization methods differ in whether or not the number of  
243 intervals  $k$  is chosen beforehand and in how the boundaries, or cut points,  
244 for the intervals are determined.

245 *Equal Frequency (EF) and Equal Width (EW) discretisation.* The *Equal Fre-*  
246 *quency* and *Equal Width* methods are probably the simplest discretization  
247 methods in use (Liu et al., 2002). With the *Equal Frequency* method, each  
248 constructed interval includes essentially the same number of data points.  
249 With the *Equal Width* method, all constructed intervals have equal length;  
250 these intervals may thus have varying numbers of data points. With both  
251 *Equal Frequency* and *Equal Width*, the parameter  $k$  dictating the number of  
252 intervals used for the discretization, is chosen beforehand. Upon discretizing  
253 the continuous variables underlying a data set, in essence different  $k$ ’s may  
254 be chosen per variable. In most applications, however, a single  $k$  is used for  
255 all variables concerned. In environmental sciences, the number of intervals is  
256 often chosen based upon expert knowledge (Chen and Pollino, 2012); without  
257 such knowledge, an appropriate number may be found by experimentation.  
258 Alternatively, the number of intervals  $k$  may be decided upon by the *Propor-*  
259 *tional k-interval Discretization* (PKID) guideline (Yang and Webb, 2009),  
260 which takes  $k = \sqrt{N}$  where  $N$  is the number of data points available.

261 The *Equal Frequency* and *Equal Width* methods are the most commonly  
262 used methods for discretizing continuous variables in environmental modeling  
263 with Bayesian networks (Aguilera et al., 2011). Chen and Pollino (2012)  
264 argue however, that the *Equal Width* method is less suited for data sets that  
265 have a markedly uneven distribution or include prominent outliers, and that  
266 the *Equal Frequency* method is less appropriate for data sets in which specific  
267 values are overrepresented.

268 *Chi-Merge (ChiM) discretization.* *Chi-Merge* is a supervised discretization  
 269 method which takes the classes associated with the available data points into  
 270 account (Kerber, 1992). The method starts by constructing a sequence of  
 271 intervals such that each interval includes a single data point. The  $\chi^2$ -statistic  
 272 is then used to decide whether two adjacent intervals be merged. For this  
 273 purpose, for each pair of adjacent intervals, the  $\chi^2$ -value is calculated from:

$$\chi^2 = \sum_{i=1}^2 \sum_{j=1}^m \frac{(A_{ij} - E_{ij})^2}{E_{ij}} \quad (1)$$

274 where  $m$  is the number of distinct classes,  $A_{ij}$  is the number of data points  
 275 in interval  $i$  that are of class  $j$ , and  $E_{ij}$  is the expected number of data points  
 276 of class  $j$  in interval  $i$  under the assumption that the class frequencies per  
 277 interval are the same; this expected number  $E_{ij}$  is established from:

$$E_{ij} = \frac{(\sum_j A_{ij}) \cdot (\sum_i A_{ij})}{\sum_i \sum_j A_{ij}} \quad (2)$$

278 In each iteration of the *Chi-Merge* method, the pair of adjacent intervals  
 279 with the smallest  $\chi^2$ -value are merged, provided that this value falls below  
 280 the confidence threshold  $\chi_{df,\alpha}^2$  read from the  $\chi^2$ -distribution table, where  $df$  is  
 281 the number of degrees of freedom  $m - 1$  and  $\alpha$  is a user-specified significance  
 282 level (preferably between 0.9 and 0.99). The iterative procedure halts when  
 283 all  $\chi^2$ -values are above the confidence threshold.

284 Since the *Chi-Merge* method serves to discretize each continuous vari-  
 285 able in a data set independently, the number of intervals constructed per  
 286 variable may differ. In order to avoid large numbers of intervals in practice,  
 287 a maximum number of intervals can be pre-set for application of the *Chi-*  
 288 *Merge* method. The iterative procedure described above is then halted as  
 289 this number of intervals is reached.

290 *Minimum Description Length Principle (MDLP) discretization.* Similar to  
 291 *Chi-Merge*, the *MDLP* method, first introduced by Fayyad and Irani (1993,  
 292 1996), is a supervised discretization method which takes the classes asso-  
 293 ciated with the data points into account. Starting with a single interval  
 294 composed of all data points sorted in increasing order of their value for the

295 variable to be discretized, *MDLP* constructs, by an iterative procedure, a se-  
 296 quence of intervals over the variable’s overall value range. Within an interval  
 297  $S$ , potential cut points  $t_i$  are defined between each pair of data values; such a  
 298 cut point would in essence partition the interval  $S$  into the two adjacent in-  
 299 tervals  $S_1^i$  and  $S_2^i$ . For each potential cut point  $t_i$  in  $S$ , the *Class Information*  
 300 *Entropy* (CIE) of the partition induced by  $t_i$  is then computed, from:

$$CIE(S, t_i) = \frac{N(S_1^i)}{N(S)} \cdot E(S_1^i) + \frac{N(S_2^i)}{N(S)} \cdot E(S_2^i) \quad (3)$$

301 where  $S_1^i, S_2^i$  are the two (sub)intervals that would be induced by the cut  
 302 point  $t_i$ ,  $N(\cdot)$  denotes the number of data points in the indicated interval, and  
 303  $E(\cdot)$  is the entropy of the class distribution estimated from the data points  
 304 in the indicated interval. This entropy  $E(S')$  for an interval  $S'$  is calculated  
 305 as:

$$E(S') = - \sum_{c_j} P_{S'}(c_j) \cdot \log_2 P_{S'}(c_j) \quad (4)$$

306 where  $c_j$  is a class and  $P_{S'}(c_j)$  is the estimated probability of occurrence  
 307 of  $c_j$  in the interval  $S'$ . The potential cut point  $t_i$  with the smallest *Class*  
 308 *Information Entropy* now is accepted as an actual cut point, provided that  
 309 doing so yields an information gain  $E(S) - CIE(S, t_i)$  satisfying the following  
 310 criterion:

$$E(S) - CIE(S, t_i) > \frac{1}{N(S)} \cdot \left( \log_2(N(S) - 1) + \Delta(S, t_i) \right) \quad (5)$$

311 where  $\Delta(S, t_i)$  equals

$$\Delta(S, t_i) = \log_2(3^m - 2) - [m \cdot E(S) - m_1 \cdot E(S_1^i) - m_2 \cdot E(S_2^i)] \quad (6)$$

312 with  $m, m_j, j = 1, 2$ , being the number of distinct classes in the intervals  
 313  $S, S_j^i$ , respectively. This procedure is repeated iteratively, with the two inter-  
 314 vals  $S_1^i$  and  $S_2^i$  substituted for the interval  $S$  for each accepted cut point  $t_i$ , as  
 315 long as there is at least one potential cut point that satisfies the information-  
 316 gain criterion above. After the procedure has halted, the accepted cut points  
 317 serve to define the intervals from the overall value range of the variable being  
 318 discretized.

319 *MDLP* discretization is one of the more commonly used discretization  
 320 methods in general (García et al., 2013). Experiments by Liu et al. (2002)

321 suggest that *MDLP* in fact is one of the best performing discretization meth-  
322 ods in practice. Despite its reported good performance, however, *MDLP* has  
323 hardly been used in environmental modeling with Bayesian networks (Fer-  
324 nandes et al., 2013).

### 325 2.3.2. The discretized data sets for the study

326 The continuous variables of the species data sets described in Section 2.1,  
327 were discretized by the four methods reviewed above, as available from the  
328 Discretization Package<sup>4</sup> of the R statistical computing software. With each  
329 of the *Equal Frequency* and *Equal Width* methods, four different discretiza-  
330 tions were constructed for the variables under study, with 3, 5, 10 and 32  
331 intervals, respectively, where the PKID criterium gave rise to the number of  
332 intervals  $k = \sqrt{N} = \sqrt{989} = 32$ . For application of the *Chi-Merge* method,  
333 a significance level of 0.99 was used as suggested in the literature; no limit  
334 was set on the number of intervals.

335 Discretization resulted in 10 data sets per species: four data sets resulted  
336 from using the *Equal Frequency* method with 3, 5, 10 and 32 intervals, re-  
337 spectively, and four resulted from using the *Equal Width* method with the  
338 same numbers of intervals; one data set resulted from application of the  
339 *Chi-Merge* method, and one set of discretized data was constructed using  
340 *MDLP*. As per species moreover, the original continuous data were used in  
341 the investigations, our study involved a total of 33 data sets.

### 342 2.4. Model learning and validation

343 From each of the 30 discretized data sets, discrete Naive Bayesian and  
344 TAN classifiers were learned. From the three continuous data sets, we con-  
345 structed NB and TAN classifiers with mixtures of truncated exponentials  
346 for the local probability distributions; the number of exponentials was set  
347 to three based on preliminary experimentation. All classifiers were learned  
348 using the Elvira software<sup>5</sup> (Elvira-Consortium, 2002). In addition, 33 logistic-  
349 regression models were constructed using the R statistical software package.

350 To arrive at reliable estimates of the predictive performance per model,  
351 a *ten-fold cross validation* procedure was used. To this end, each data set  
352  $D$  was partitioned into ten equally-sized disjoint subsets, or *folds*,  $D_i$ ,  $i =$   
353  $1, \dots, 10$ . Then, for each fold, the following procedure was run:

---

<sup>4</sup><https://cran.r-project.org/web/packages/discretization/index.html>

<sup>5</sup>Elvira is a public-domain Java-based software package: <http://leo.ugr.es/elvira>

- 354 • set the current fold  $D_i$  aside for testing;
- 355 • learn the appropriate type of classification model from the set  $D^{-i} =$   
356  $\bigcup_{j=1, \dots, 10, j \neq i} D_j$  composed of the data from the other nine folds;
- 357 • estimate the performance of the thus learned model by classifying the  
358 data points from  $D_i$ .

359 The predictive performance of the classifier learned from the entire data set  
360  $D$  now is estimated as the performance result averaged over the ten runs.

361 As measures of performance for the learned classification models, the  
362 well-known *sensitivity* and *specificity* characteristics are used. Estimates for  
363 these characteristics are calculated from:

$$\text{sensitivity} = \frac{TP}{TP + FN} \quad (7)$$

364

$$\text{specificity} = \frac{TN}{FP + TN} \quad (8)$$

365 where  $TP$  is the number of *true positives*, that is, the number of data points  
366 in the test set in which the species is known to be *present* and which are  
367 assigned to the *present* class by the learned classifier, and  $TN$  is the num-  
368 ber of *true negatives*, that is, the number of data points in the test set in  
369 which the species is *absent* and which are assigned to the *absent* class by the  
370 classifier;  $FP$  is the number of *false positives*, that is, the number of data  
371 points in the test set which are classified as *present*, yet are known to be  
372 *absent*, and  $FN$  is the number of *false negatives*, that is, the number of data  
373 points which are classified as *absent* and are known to be *present*. The thus  
374 obtained *sensitivity* and *specificity* estimates are combined into an *averaged*  
375 *performance* estimate through

$$\text{averaged performance} = \frac{1}{2} (\text{sensitivity} + \text{specificity}) \quad (9)$$

376 Since the *sensitivity* and *specificity* estimates found for a classifier are  
377 dependent of the decision threshold used for classification, all Bayesian-  
378 network classifiers were validated using maximum-probability classification to  
379 allow a fair comparison with their matching logistic-regression models. The  
380 Bayesian-network classifiers were also validated using probability-threshold  
381 classification with the prevalences of the various bird species as the decision  
382 thresholds. All in all, with each of the 33 data sets, five classification models

Table 2: Minimum, maximum and mean number of intervals constructed by the *Chi-Merge* and *MDLP* discretization methods, per species data set.

		<i>Cecropis daurica</i>	<i>Turdus viscivorus</i>	<i>Accipiter nisus</i>
<i>Chi-Merge</i>	Mean	52.9	12	18.7
	Minimum	14	6	9
	Maximum	91	25	29
<i>MDLP</i>	Mean	1.5	2.4	1.8
	Minimum	1	1	1
	Maximum	2	4	3

383 were learned and validated: an NB, a TAN and a logistic-regression model  
 384 were constructed and evaluated using maximum-probability classification,  
 385 and an NB and a TAN were learned and validated using probability-threshold  
 386 classification.

### 387 3. Results

388 The experimental results from using the different discretization methods  
 389 on the various data sets are summarized by the granularity of the resulting  
 390 discretizations and by the performance of the learned classification models.

#### 391 3.1. Granularity of discretization

392 For use of the *Equal Frequency* and *Equal Width* discretization methods,  
 393 the numbers of intervals to be constructed for a continuous variable were  
 394 chosen beforehand, as 3, 5, 10 and 32, respectively; for each variable, the same  
 395 number of intervals was used. With the *Chi-Merge* and *MDLP* methods, the  
 396 numbers of intervals to be constructed were not pre-set but rather established  
 397 by the methods themselves, for each variable separately. Table 2 reports,  
 398 for each species data set, the numbers of intervals constructed by the *Chi-*  
 399 *Merge* and *MDLP* methods, respectively; the means reported in the table  
 400 were calculated by averaging over the discretizations of all variables in the  
 401 data set at hand. The table shows that, for our data sets, the *MDLP* method  
 402 resulted in quite coarse discretizations, with just a limited number of intervals  
 403 per variable. The *Chi-Merge* method, on the other hand, resulted in more  
 404 fine-grained discretizations, with over 50 intervals for some of the variables  
 405 involved.

406 *3.2. Maximum-probability classification*

407 From each of the species data sets, Bayesian-network classifiers and logistic-  
408 regression models were learned as described in Section 2.4. The performances  
409 of the learned models using maximum-probability classification are visualized  
410 in Figure 3, which shows the *sensitivity* and *specificity* estimates found; Ta-  
411 ble 3 summarizes these estimates in the models' *averaged performances*.

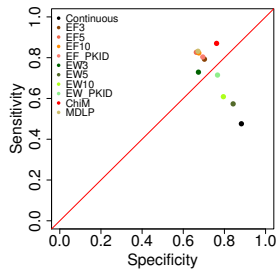
412 For the data set pertaining to *Turdus viscivorus*, all logistic-regression  
413 models showed quite similar performance, regardless of whether continuous  
414 or discretized data were used and, in the latter case, regardless of the dis-  
415 cretization method employed (Fig. 3(c)). The Bayesian-network classifiers  
416 (Figs. 3(a) and 3(b)) showed more divergence in their performance charac-  
417 teristics. From among the discretization methods used, the *Chi-Merge* method  
418 resulted in the best balance of the *specificity* and *sensitivity* characteris-  
419 tics estimated for the classifiers, with an *averaged performance* of 0.82. Figs. 3(a)  
420 and 3(b) further show that the continuous Bayesian-network classifiers had  
421 a worse *sensitivity* than the classifiers learned from discretized data.

422 While the data set pertaining to *Turdus viscivorus* is well balanced with  
423 respect to the two classes, the other two data sets are less balanced, with a  
424 prevalence of 84% for the *Cecropis daurica* and a prevalence of 27% for the  
425 *Accipiter nisus*, respectively. For these less balanced data sets, all constructed  
426 models were found to excel at predicting the most probable class. More  
427 specifically, for the *Cecropis daurica* all classifiers attained a high *sensitivity*  
428 (Figs. 3(d), 3(e) and 3(f)), while for the *Accipiter nisus* the classifiers attained  
429 a high *specificity* (Figs. 3(g), 3(h) and 3(i)).

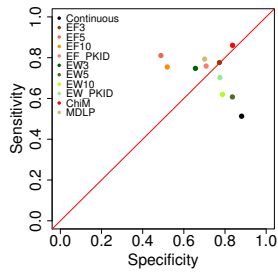
430 For the data set pertaining to *Cecropis daurica*, all Bayesian-network  
431 classifiers showed quite similar performance, yet with a notable single excep-  
432 tion. The Naive Bayesian classifier learned from the data after discretization  
433 with the *Chi-Merge* method, showed very good performance in terms of both  
434 *sensitivity* and *specificity*; this classifier in fact resulted in an *averaged perfor-*  
435 *mance* of 0.93 (Table 3). Also for the *Accipiter nisus* data set discretized with  
436 the *Chi-Merge* method, did the NB classifier show the best balance of the  
437 *sensitivity* and *specificity* estimates attained, with an *averaged performance*  
438 of 0.84. While the TAN classifier never reached a *sensitivity* higher than 0.6  
439 for the *Accipiter nisus* data set, the NB classifier gave *sensitivity* estimates  
440 higher than 0.7 after discretizing the data with the *Equal Frequency* method  
441 with 3, 5 and 10 intervals, with *MDLP* and with the *Chi-Merge* method.

442 The performance characteristics of the logistic-regression models learned  
443 from the *Cecropis daurica* and *Accipiter nisus* data sets (Figs. 3(f) and 3(i))

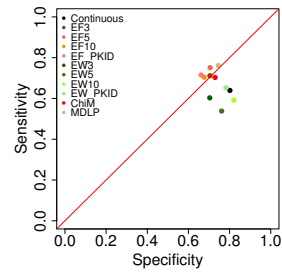




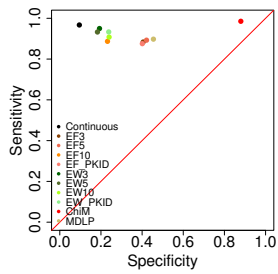
(a) *Turdus viscivorus* NB



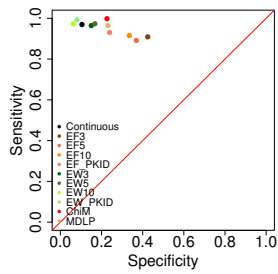
(b) *Turdus viscivorus* TAN



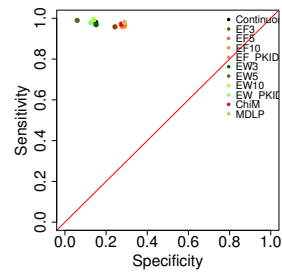
(c) *Turdus viscivorus* LR



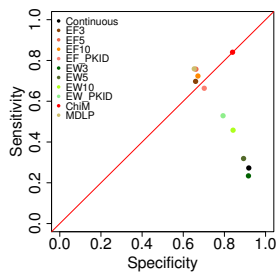
(d) *Cecropis daurica* NB



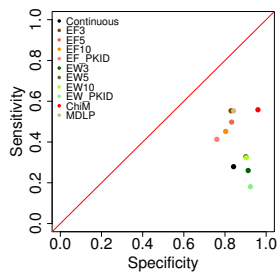
(e) *Cecropis daurica* TAN



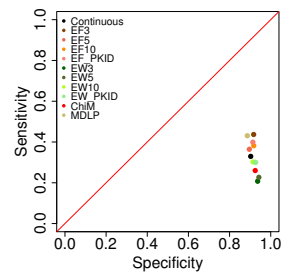
(f) *Cecropis daurica* LR



(g) *Accipiter nisus* NB



(h) *Accipiter nisus* TAN



(i) *Accipiter nisus* LR

Figure 3: *Sensitivity* and *specificity* estimates for the NB, TAN and logistic-regression (LR) models with maximum-probability classification, per species data set.

Table 3: Averaged performance estimates of the classification models with maximum-probability classification, per species data set.

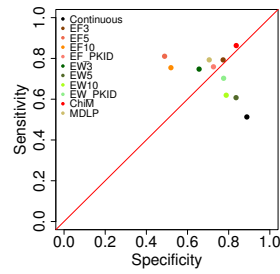
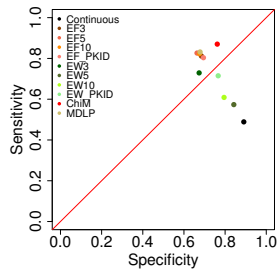
	<i>Turdus viscivorus</i>			<i>Cecropis daurica</i>			<i>Accipiter nisus</i>		
	NB	TAN	LR	NB	TAN	LR	NB	TAN	LR
Continuous	0.68	0.69	0.72	0.53	0.54	0.57	0.60	0.54	0.59
EF3	0.75	0.75	0.71	0.64	0.67	0.60	0.68	0.69	0.68
EF5	0.75	0.74	0.73	0.66	0.64	0.63	0.71	0.67	0.63
EF10	0.75	0.75	0.69	0.66	0.63	0.62	0.70	0.63	0.65
EF_PKID	0.75	0.68	0.69	0.64	0.59	0.62	0.68	0.59	0.66
EW3	0.70	0.71	0.65	0.57	0.56	0.57	0.58	0.59	0.57
EW5	0.70	0.73	0.65	0.56	0.57	0.53	0.61	0.62	0.58
EW10	0.70	0.70	0.71	0.58	0.52	0.57	0.65	0.61	0.61
EW_PKID	0.74	0.74	0.72	0.59	0.54	0.55	0.66	0.55	0.62
ChiM	0.82	0.82	0.72	0.93	0.61	0.63	0.84	0.76	0.59
MDLP	0.75	0.78	0.75	0.68	0.60	0.63	0.71	0.70	0.66

444 again were hardly affected by the discretization method used. For these  
 445 two data sets, the *averaged performance* estimates found with the logistic-  
 446 regression models were in the 0.53 – 0.68 range (Table 3). While for all  
 447 models very good performance at predicting the most probable class was  
 448 seen, the best *specificity* achieved by these models for *Cecropis daurica* was  
 449 smaller than 0.3; similarly, the best *sensitivity* achieved for *Accipiter nisus*  
 450 by the logistic-regression models was below 0.45.

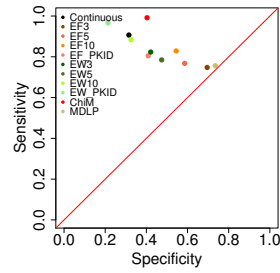
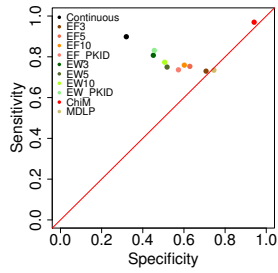
### 451 3.3. Probability-threshold classification

452 The performances of the Bayesian-network classifiers learned from each of  
 453 the species data sets are once more investigated, this time using probability-  
 454 threshold classification. The detailed results are visualized in Fig. 4, in terms  
 455 of the *sensitivity* and *specificity* estimates found; Table 4 summarizes these  
 456 estimates in the models’ *averaged performances*.

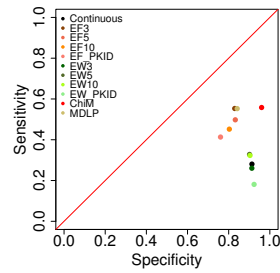
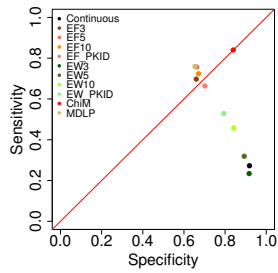
457 For the data set pertaining to *Turdus viscivorus*, the learned Bayesian-  
 458 network classifiers were found to exhibit similar performance with probability-  
 459 threshold classification as with maximum-probability classification (Figs. 4(a)  
 460 and 4(b)). Since the *Turdus viscivorus* data set includes a binary class  
 461 variable, maximum-probability classification was equivalent to probability-  
 462 threshold classification with a decision threshold equal to 0.5. Based on  
 463 the prevalence for *Turdus viscivorus*, probability-threshold classification was  
 464 performed with a threshold of 0.47. Given the small difference in decision



(a) *Turdus viscivorus* NB (b) *Turdus viscivorus* TAN



(c) *Cecropis daurica* NB (d) *Cecropis daurica* TAN



(e) *Accipiter nisus* NB (f) *Accipiter nisus* TAN

Figure 4: *Sensitivity* and *specificity* estimates for the NB and TAN classifiers with probability-threshold classification, per species data set.

Table 4: Averaged performance estimates of the classifiers using probability-threshold classification, per species data set.

	<i>Turdus viscivorus</i>		<i>Cecropis daurica</i>		<i>Accipiter nisus</i>	
	NB	TAN	NB	TAN	NB	TAN
Continuous	0.68	0.69	0.61	0.61	0.60	0.54
EF3	0.75	0.75	0.72	0.72	0.68	0.69
EF5	0.75	0.74	0.73	0.69	0.71	0.67
EF10	0.75	0.75	0.73	0.69	0.70	0.63
EF_PKID	0.75	0.68	0.65	0.61	0.68	0.59
EW3	0.70	0.71	0.63	0.62	0.58	0.59
EW5	0.70	0.73	0.63	0.63	0.61	0.62
EW10	0.70	0.70	0.64	0.61	0.65	0.61
EW_PKID	0.74	0.74	0.65	0.59	0.66	0.55
ChiM	0.82	0.82	0.96	0.70	0.84	0.76
MDLP	0.75	0.78	0.74	0.75	0.71	0.70

465 threshold used, similar performance of the Bayesian-network classifiers under  
 466 the two types of classification was not unexpected.

467 Also for the data set pertaining to *Accipiter nisus* were the performances  
 468 of the Bayesian-network classifiers with probability-threshold classification  
 469 comparable to those found with maximum-probability classification (Figs. 4(e)  
 470 and 4(f)). With this data set, however, the decision threshold for classifi-  
 471 cation was set to the prevalence of 0.27 of the bird species, which differed  
 472 substantially from the 0.5 threshold used with maximum-probability classi-  
 473 fication.

474 When validated on the data set pertaining to *Cecropis daurica*, the Bayesian-  
 475 network classifiers showed a different performance with probability-threshold  
 476 classification (Figs. 4(c) and 4(d)) than with maximum-probability classifica-  
 477 tion. In fact, use of the species' prevalence of 0.84 for the decision threshold  
 478 for classification resulted in a better balance of the *sensitivity* and *speci-*  
 479 *ficity* characteristics estimated for the classifiers (Table 4) than use of the 0.5  
 480 threshold with maximum-probability classification. For the Naive Bayesian  
 481 classifiers specifically, the good performances in terms of *sensitivity* were  
 482 matched by a *specificity* between 0.7 and 0.8 after discretizing the data with  
 483 *MDLP* and with the *Equal Frequency* method with three intervals; the corre-  
 484 sponding *specificity* estimates found with maximum-probability classification  
 485 were below 0.4. For the NB classifier moreover, discretization of the data with  
 486 the *Chi-Merge* method gave the best *averaged performance* estimate, equal

487 to 0.96. For the TAN classifier, discretization with the *MDLP* method gave  
488 the best overall result.

489 From the detailed *sensitivity* and *specificity* estimates plotted for the var-  
490 ious Bayesian-network classifiers in Figs. 3 and 4, a general pattern emerges.  
491 Both with maximum-probability classification and with probability-threshold  
492 classification, discretization of the data with the *Equal Width* method tends  
493 to result in classifiers with a good performance at predicting the most prob-  
494 able class, that is, the *C. daurica* being *present* and the *A. nisus* being  
495 *absent*. With both types of classification, moreover, discretization with the  
496 *Equal Frequency* method tends to result in a better balance of the *sensitivity*  
497 and *specificity* characteristics of the learned Bayesian-network classifiers.

#### 498 4. Discussion

499 Based on the experimental results described in Section 3, we discuss some  
500 implications for use of the various discretization methods and classification  
501 models in species distribution modeling.

502 *Logistic-regression models.* Regression methods are widely used in environ-  
503 mental modeling in general (Schmitz et al., 2005) and for species distribution  
504 modeling in particular (Li and Wang, 2013). For well-balanced data sets, in  
505 which a species is (more or less) equally likely to be *present* as it is to be  
506 *absent*, our experimental results suggest that logistic-regression models can  
507 attain relatively high *sensitivity* and *specificity* characteristics. The overall  
508 performance of these models moreover, appears not to be affected by dis-  
509 cretization of the data nor by the method used if the data were discretized.  
510 For less balanced data sets, however, logistic-regression models tend to fail  
511 at predicting the least probable class.

512 From an environmental point of view, a species distribution model should  
513 accurately predict the presence of a specific species in a territory, that is, it  
514 should show a high *sensitivity*. For abundant species, such as *Cecropis dau-*  
515 *rica* in our study, logistic-regression models can indeed attain a high *sensi-*  
516 *tivity* and thereby show satisfactory performance. In real-world applications  
517 however, attention will mostly focus on endangered or rare species, such  
518 as *Accipiter nisus*. Our experimental results suggest that, for such species,  
519 logistic-regression models may not be able to achieve a satisfactory perfor-  
520 mance. For such species, therefore, using logistic regression may not be the  
521 best possible choice.

522 *Bayesian-network classifiers and the effect of decision thresholds.* One of the  
523 advantages of Bayesian-network classifiers over other types of classifier is that  
524 the classification decision is separated from the prediction process (Uusitalo,  
525 2007). The Bayesian networks underlying these classifiers in essence return  
526 a posterior probability distribution over the class variable given the case ob-  
527 servations, based upon which a classification decision is taken. As discussed  
528 in the previous sections, cases can then be assigned to the most probable  
529 class or to a class decided upon through a probability threshold.

530 Naive Bayesian and TAN classifiers are known to show a tendency to pro-  
531 duce rather skewed posterior distributions for their class variable (Bennett,  
532 2000). As a consequence, choosing non-extreme thresholds with probability-  
533 threshold classification may not dramatically change performance compared  
534 to maximum-probability classification. For the *Turdus viscivorus* and *Ac-  
535 cipiter nisus* data sets in our study, in fact, classification with the decision  
536 thresholds of 0.47 and 0.27, respectively, did not result in a performance  
537 different from using the 0.5 threshold of maximum-probability classification.  
538 For the former species, this experimental finding was not unexpected given  
539 the small difference between the thresholds of 0.47 and 0.5. For the latter  
540 species, the difference between the two thresholds involved was more substan-  
541 tial. The finding of similar performance with the two types of classification  
542 now indicates that, for none or just a few cases, the established posterior  
543 probability of the species being *present* was in the 0.27 – 0.5 range, which  
544 would indeed be explained by the tendency of the Bayesian-network classi-  
545 fiers to produce rather skewed distributions over their class variable. While  
546 for *Turdus viscivorus* and *Accipiter nisus* using the species' prevalence for the  
547 decision threshold did not have any impact on classification performance, for  
548 the *Cecropis daurica* species the performance of the Bayesian-network clas-  
549 sifiers did improve with probability-threshold classification using the more  
550 extreme decision threshold probability of 0.84.

551 The above insights from our experimental results suggest that using pro-  
552 bability-threshold classification can be beneficial with Bayesian-network clas-  
553 sifiers developed for species with quite small prevalences.

554 *Continuous Bayesian-network classifiers.* Direct use of available continuous  
555 data is often recommended for Bayesian-network learning, to avoid loss of  
556 information due to discretization (Uusitalo, 2007). The current generation of  
557 Bayesian networks can cope with continuous probability distributions only  
558 to some extent, however: local distributions for the continuous variables are

559 required to be Gaussian (Lauritzen and Wermuth, 1989) or are approximated  
560 by polynomial or exponential functions, such as the MTEs used in our study.

561 In our experimental study, the Bayesian-network classifiers with MTEs  
562 for their local distributions showed good performance at predicting the most  
563 probable class. Since typically a large number of data points is required to  
564 allow satisfactory approximation of the continuous distributions at hand, this  
565 good performance may be attributed, to at least some extent, to the availabil-  
566 ity of many data points from the predominant class. For endangered or rare  
567 species, where the class of interest is the less probable one, our experimental  
568 results suggest that Bayesian-network classifiers with MTEs may result in  
569 relatively poor *sensitivity* and hence show unsatisfactory performance. For  
570 such species, direct use of available continuous data may not be the best  
571 choice for finding Bayesian-network classifiers of good performance.

572 *Unsupervised discretization.* The unsupervised *Equal Width* and *Equal Fre-*  
573 *quency* discretization methods are widely used in environmental modeling  
574 through Bayesian networks (Aguilera et al., 2011). Chen and Pollino (2012)  
575 already argued that both methods are suitable for discretizing variables with  
576 a more or less even distribution over their values. They further argued that  
577 use of the *Equal Width* method is less appropriate for data sets that have  
578 a markedly uneven distribution or include prominent outliers, and that the  
579 *Equal Frequency* method is less suited for data sets in which specific values  
580 are overrepresented. The land-use variables in our study typically do not  
581 have even distributions, as was illustrated for the *Olive cropland* variable in  
582 Fig. 2(b).

583 The *Equal Frequency* method partitions the overall value range of a con-  
584 tinuous variable into  $k$  intervals such that each interval includes an essentially  
585 equal number of data points. Yet, data points with the same value for the  
586 continuous variable to be discretized are never placed in different intervals.  
587 Since the feature variables in our study capture types of land use that are  
588 present in relatively few grid cells, for any such variable a large number of  
589 data points include the value 0. These data points are all included in the  
590 first interval therefore, and the remaining data points are equally distributed  
591 over the remaining  $k - 1$  intervals.

592 Our experimental results indicate that, with *Equal Frequency* discretiza-  
593 tion, all constructed Bayesian-network classifiers show good performance at  
594 predicting the most probable class; this good performance is generally bal-  
595 anced by a reasonable performance for the less probable class. The results

596 further show that using just three intervals for the discretization tends to  
597 result in the best balance of the *sensitivity* and *specificity* characteristics for  
598 all Bayesian-network classifiers learned. Since for most land-use variables a  
599 large number of data points include the value 0, the majority of the intervals  
600 constructed with the *Equal Frequency* method will include just a few data  
601 points. In fact, the more intervals are used, the fewer data points are ex-  
602 pected per interval and the less informative the intervals tend to become for  
603 classification purposes. Using a small number of intervals therefore appears  
604 to be the best option upon *Equal Frequency* discretization of data sets in  
605 which specific values are overrepresented.

606 The *Equal Width* discretization method partitions the overall value range  
607 of a continuous variable into  $k$  intervals such that all intervals are of equal  
608 length. Just like the *Equal Frequency* method, it includes all data points with  
609 the value 0 for the variable to be discretized in the first interval. Since for  
610 most land-use variables a large number of data points include this value and  
611 the method further aims at constructing intervals of equal length, actually  
612 the majority of data points will be included in this first interval and very  
613 few points remain for the subsequent intervals, which causes these intervals  
614 to be rather uninformative.

615 Our experimental results now indicate that, with *Equal Width* discretiza-  
616 tion, all constructed Bayesian-network classifiers show good performance at  
617 predicting the most probable class, just as with *Equal Frequency* discretiza-  
618 tion. With *Equal Width* discretization however, this good performance is  
619 balanced by a relatively poor performance for the less probable class, as a  
620 consequence of the constructed highly dominant first interval. While, with  
621 *Equal Frequency* discretization, using three intervals resulted in the best bal-  
622 ance of the *sensitivity* and *specificity* characteristics for all Bayesian-network  
623 classifiers learned, the experimental results obtained with *Equal Width* dis-  
624 cretization suggest that using a small number of intervals may not always  
625 give a well-balanced performance. For the *Accipiter nisus* data set, with the  
626 low prevalence of its species, in fact, using three intervals for the discretiza-  
627 tion resulted in a very high specificity while more intervals were required to  
628 attain a reasonable sensitivity.

629 *Supervised discretization.* The supervised *Chi-Merge* and *MDLP* methods  
630 take the classes associated with the available data points into account upon  
631 discretizing the continuous variables involved. The two methods differ in  
632 their starting points for the iterative procedure and in their criteria for merg-



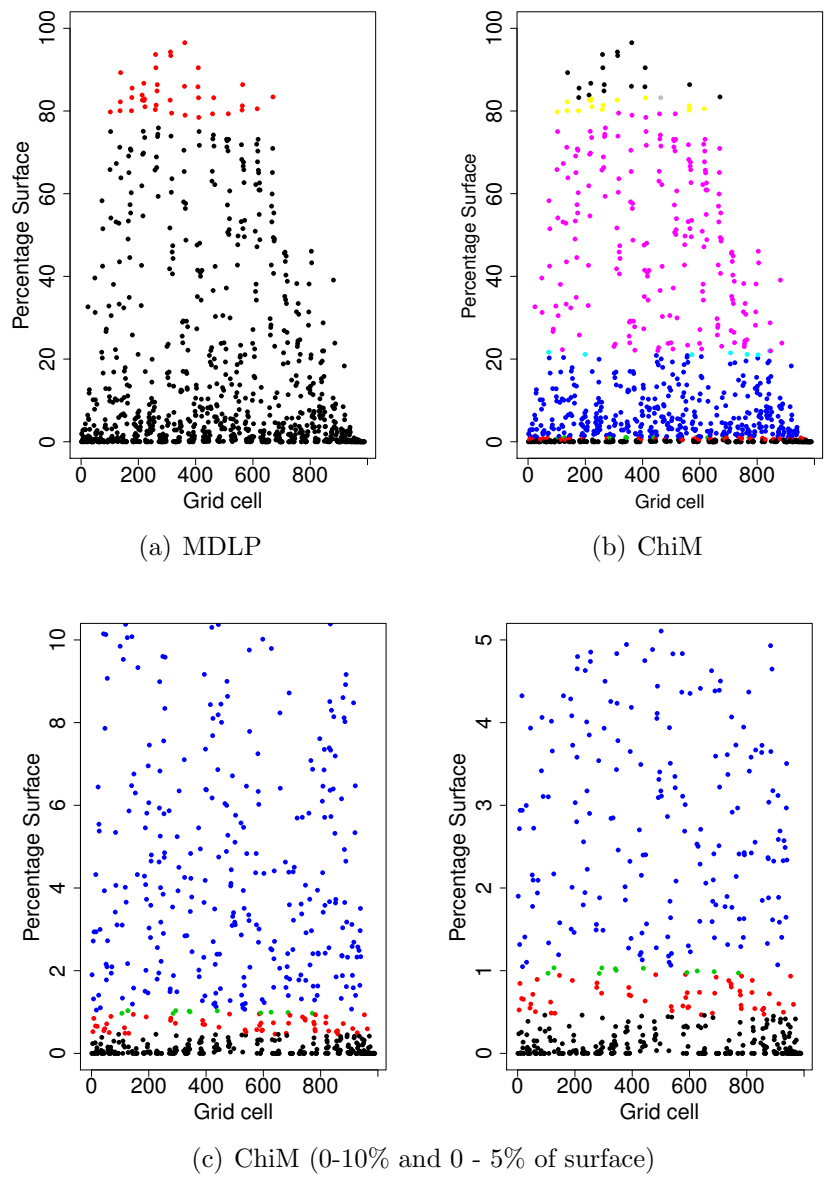


Figure 5: Distribution of the *Olive croplands* variable, discretized by the *MDLP* method (a) and by the *Chi-Merge* method (b), with a detailed view of the latter discretization for 0-10% and 0-5%, respectively, of the surface (c).

633 ing and splitting intervals. The *Chi-Merge* method starts with a separate  
634 interval per data point and iteratively merges two adjacent intervals if the  
635 class distributions in these intervals are more or less similar. The *MDLP*  
636 method on the other hand, starts with a single interval including all data  
637 points and iteratively splits an interval if the class distributions in the re-  
638 sulting subintervals are more skewed than the distribution in the original  
639 interval.

640 As argued above, the continuous land-use variables in our study have  
641 highly skewed distributions, as a result of the heterogeneous conditions of  
642 Andalusia. From among the two supervised discretization methods reviewed  
643 in our study, the *Chi-Merge* method seems better able to capture the charac-  
644 teristics of the data than the *MDLP* method. As an example, Fig. 5 depicts  
645 the available data points in terms of their *Olive croplands* coverage, for the  
646 cells of the UTM grid. The full range of the percentage of surface covered,  
647 is partitioned into intervals by the two discretization methods. The dis-  
648 cretization constructed by the *MDLP* method is shown in Fig. 5(a) and the  
649 discretization found by *Chi-Merge* is shown in Fig. 5(b); the various parti-  
650 tions are indicated in color. Fig. 5(a) reveals that, with the *MDLP* method,  
651 just two intervals were constructed for the entire range of the percentage of  
652 covered surface. With the *Chi-Merge* method, multiple intervals were cre-  
653 ated: four intervals were constructed for the lower percentages of surface  
654 coverage (Fig. 5(c)) and four more intervals resulted for the remainder of the  
655 percentage range. The difference between the resulting discretizations may,  
656 to some extent, be due to the stopping criteria employed by the two meth-  
657 ods. Yet also the tendency of the *MDLP* method to construct intervals with  
658 class distributions of low entropy may cause this method to be less sensi-  
659 tive to small shifts in already quite skewed distributions than the *Chi-Merge*  
660 method is.

661 For all species data sets discretized with the *Chi-Merge* method, the Naive  
662 Bayesian classifiers attained high *sensitivity* and *specificity* characteristics,  
663 with averaged performances between 0.82 and 0.96. A similar trend was  
664 seen for the TAN classifiers constructed from the *Turdus viscivorus* data set  
665 discretized with *Chi-Merge*. For the less balanced data sets discretized with  
666 the *Chi-Merge* method, the TAN classifiers excelled at predicting the most  
667 probable class. For the least probable class, however, TAN classifiers with a  
668 better performance resulted from discretizing the data with less sophisticated  
669 methods. For the *Cecropis daurica* data set in fact, using the *Chi-Merge*  
670 method for discretization resulted in TANs with averaged performances of

671 0.61 and 0.70, while the best performing TANs had averaged performances  
672 of 0.67 and 0.72, respectively. The lesser performance found from using the  
673 *Chi-Merge* method with the *Cecropis daurica* data set may be attributed to  
674 the relatively large number of intervals constructed for the various feature  
675 variables: some of these intervals are likely to include only very few data  
676 points and, as a consequence, the strengths estimated for the dependencies  
677 involved in the TANs will most likely be unreliable.

678 For the species data sets discretized with the *MDLP* method, the per-  
679 formance trends of all Bayesian-network classifiers were more or less similar  
680 to those found for the sets discretized with the *Chi-Merge* method, although  
681 less prominent. Overall, the averaged performances of the various classifiers  
682 were found to lie below those of the corresponding classifiers learned from  
683 the *Chi-Merge* discretized data.

## 684 5. Conclusions and future research

685 In our experimental study, we compared the performances of different  
686 types of classification model and different discretization methods in view of  
687 species distribution modeling. In the study, we focused on prediction of the  
688 presence of various bird species in Andalusia from land-use data, and con-  
689 sidered to this end three species with different prevalence rates. The experi-  
690 mental results obtained suggest that Bayesian-network classifiers, and among  
691 these especially the Naive Bayesian classifiers, may be preferable to logistic-  
692 regression models for the environmental-science context at hand. Our results  
693 further indicate that the *Chi-Merge* method may be the preferred method  
694 for discretizing the continuous variables involved, since with this method the  
695 best averaged performance results in terms of both *sensitivity* and *specificity*  
696 were found. As it is a supervised method, it is computationally more in-  
697 volved than the better known *Equal Frequency* and *Equal Width* methods  
698 for discretization. Implementations of the *Chi-Merge* method are available  
699 in software packages such as R for ready use in practice.

700 While most applications of Bayesian networks require discretization of the  
701 continuous variables underlying available data, only a restricted set of meth-  
702 ods are used in practice. For species distribution modeling through Bayesian  
703 networks more specifically, further research efforts are required to gain insight  
704 in the foundational properties of the various discretization methods proposed  
705 in the literature and to establish their practical properties upon application  
706 to different types of environmental data. While the conclusions obtained from

707 our experimental study are likely to hold for data sets with similar charac-  
708 teristics as our land-use data, the results cannot be directly extrapolated  
709 to other environmental data, such water quality, air pollution and climatic  
710 data, without further study. Moreover, since expert knowledge is often taken  
711 as a primary source of information in environmental science (Henriksen et al.,  
712 2007), and is in fact used for choosing cut points for discretization, the quality  
713 of expert-based discretizations should be compared with the discretizations  
714 found with automated methods.

715 From a wider future perspective, it is worthwhile to study the strengths  
716 and weaknesses of using Bayesian networks for species distribution modeling  
717 compared to using the more common domain-specific models proposed in the  
718 literature, such as BIOCLIM and FLORAMAP.

## 719 **Acknowledgements**

720 R.F. Ropero is supported by the FPU research grant, AP2012-2117,  
721 funded by the Spanish Ministry of Education, Culture and Sport. The re-  
722 search presented in this paper was performed while R.F. Ropero stayed at  
723 Utrecht University; this research stay was supported by the Spanish Ministry  
724 of Education, Culture and Sport (grant EST15/00067).

## 725 **References**

- 726 Aguilera, P. A., Fernández, A., Fernández, R., Rumí, R., Salmerón, A., 2011.  
727 Bayesian networks in environmental modelling. *Environmental Modelling*  
728 & *Software* 26, 1376–1388.
- 729 Aguilera, P. A., Fernández, A., Ropero, R. F., Molina, L., 2013. Ground-  
730 water quality assessment using data clustering based on hybrid Bayesian  
731 networks. *Stochastic Environmental Research & Risk Assessment* 27 (2),  
732 435–447.
- 733 Baur, B., Bozdog, S., 2015. A canonical correlation analysis-based dynamic  
734 Bayesian network prior to infer gene regulatory networks from multiple  
735 types of biological data. *Journal of Computational Biology* 22, 289–299.
- 736 Bennett, P., 2000. Assessing the calibration of naive Bayes’ posterior esti-  
737 mates. Tech. Rep. CMU-CS00-155, Carnegie Mellon University.

- 738 Busby, J., 1986. Bioclimate prediction system (BIOCLIM). Users manual  
739 version 2.0. Australian Biological Resources, Study Leaflet, Canberra, Aus-  
740 tralia.
- 741 Chen, S. H., Pollino, C. A., 2012. Good practice in Bayesian network mod-  
742 elling. *Environmental Modelling & Software* 37, 134–145.
- 743 Chow, C. K., Liu, C. N., 1968. Approximating discrete probability distribu-  
744 tions with dependence trees. *IEEE Transactions on Information Theory*  
745 14, 462–467.
- 746 Davison, A. C., Ramesh, N., 1996. Some models for discretized series of  
747 events. *Journal of the American Statistical Association* 91, 601–609.
- 748 Dedecker, A. P., Goethals, P. L. M., Gabriels, W., De Pauw, N., 2004. Opti-  
749 mization of Artificial Neural Network (ANN) model design for prediction  
750 of macroinvertebrates communities in the Zwalm river basin (Flanders,  
751 Belgium). *Ecological Modelling* 174, 161–173.
- 752 Dyer, F., ElSawah, S., Croke, B., Griffiths, R., Harrison, E., Lucena-Moya,  
753 P., Jakeman, A. J., 2014. The effects of climate change on ecologically-  
754 relevant flow regime and water quality attributes. *Stochastic Environmen-  
755 tal Research & Risk Assessment* 28, 67–82.
- 756 Elith, J., Graham, C. H., Anderson, R. P., Dudik, M., Ferrier, S., Guisan,  
757 A., Hijmans, R. J., Huettmann, F., Leathwick, J. R., Li, J., Lohmann,  
758 L. G., Loiselle, B. A., Manion, G., Moritz, C., Nakamura, M., Nakazawa,  
759 Y., Overton, J. M., Peterson, A. T., Phillips, S. J., Richardson, K. S.,  
760 Scachetti-Pereria, S., Schapire, R. E., Soberón, J., Williams, S., Wisz,  
761 M. S., Zimmermann, N. E., 2006. Novel methods to improve prediction of  
762 species’ distribution from occurrence data. *Ecography* 29, 129–151.
- 763 Elvira-Consortium, 2002. Elvira: An environment for creating and using  
764 probabilistic graphical models. In: *Proceedings of the First European  
765 Workshop on Probabilistic Graphical Models*. pp. 222–230.  
766 URL <http://leo.ugr.es/elvira>
- 767 Fayyad, U., Irani, K., 1993. Multi-interval discretization of continuous-valued  
768 attributes for classification learning. In: *Thirteenth International Joint  
769 Conference on Artificial Intelligence*. San Mateo, CA: Morgan Kaufmann.  
770 pp. 1022–1027.

- 771 Fayyad, U., Irani, K., 1996. Discretizing continuous attributes while learning  
772 Bayesian networks. In: Thirteenth International Conference on Machine  
773 Learning. Morgan Kaufmann. pp. 157–165.
- 774 Fernandes, J. A., Lozano, J. A., Inza, I., Irigoien, X., Pérez, A., Rodríguez,  
775 J. D., 2013. Supervised pre-processing approaches in multiple class variables  
776 clasification for fish recruitment forecasting. *Environmental Modelling &  
777 Software* 40, 245–254.
- 778 Friedman, N., Geiger, D., Goldszmidt, M., 1997. Bayesian network classifiers.  
779 *Machine Learning* 29, 131–163.
- 780 Fukuda, S., De Baets, B., Waegeman, W., Verwaeren, J., Mouton, A.,  
781 2013. Habitat prediction and knowledge extraction for spawning Euro-  
782 pean grayling (*Thymallus thymallus* L.) using a broad range of species  
783 distribution models. *Environmental Modelling & Software* 47, 1–6.
- 784 García, S., Luengo, J., Saez, J. A., Lopez, V., Herrera, F., 2013. A survey of  
785 discretization techniques: Taxonomy and empirical analysis in supervised  
786 learning. In: *IEEE Transactions on Knowledge and Data Engineering*. pp.  
787 734–750.
- 788 Henriksen, H. J., Rasmussen, P., Brandt, G., von Bülow, D., Jensen, F. V.,  
789 2007. Public participation modelling using Bayesian networks in manage-  
790 ment of groundwater contamination. *Environmental Modelling & Software*  
791 22, 1101–1113.
- 792 Jensen, F. V., Nielsen, T. D., 2007. *Bayesian Networks and Decision Graphs*.  
793 Springer.
- 794 Jones, P., Gladkov, A., 1999. FloraMap: A computer tool for predicting the  
795 distribution of plants and other organisms in the Wild. Centro Interna-  
796 cional de Agricultura Tropical (CIAT), Cali, Colombia.
- 797 Kerber, R., 1992. Chimerge: Discretization of numeric attributes. In: *AAAI-  
798 92, Ninth National Conference Artificial Intelligence*. AAAI Press/The  
799 MIT Press. pp. 123–128.
- 800 Lachiche, N., Flach, P., 2003. Improving accuracy and cost of two-class and  
801 multi-class probabilistic classifiers using ROC curves. In: Fawcett, T.,

- 802 Mishra, N. (Eds.), Proceedings of the Twentieth International Conference  
803 on Machine Learning. AAAI Press, Menlo Park, pp. 416–423.
- 804 Langseth, H., Nielsen, T. D., Rumí, R., Salmerón, A., 2012. Mixtures of  
805 Truncated Basis Functions. *International Journal of Approximate Reasoning* 53 (2), 212–227.  
806
- 807 Lauritzen, S. L., Jensen, F., 2001. Stable local computation with conditional  
808 Gaussian distributions. *Statistics and Computing* 11, 191–203.
- 809 Lauritzen, S. L., Wermuth, N., 1989. Graphical models for associations between  
810 variables, some of which are qualitative and some quantitative. *The Annals of Statistics* 17, 31–57.  
811
- 812 Li, X., Wang, Y., 2013. Applying various algorithms for species distribution  
813 modelling. *Integrative Zoology* 8(2), 124–135.
- 814 Li, Y., 2007. Control of spatial discretisation in coastal oil spill modelling.  
815 *International Journal of Applied Earth Observation* 9, 392–402.
- 816 Lima, M. D., Nassar, S. M., Rodrigues, P. I. R., Freitas-Filho, P., Jacinto,  
817 C. M. C., 2014. Heuristic discretization method for Bayesian networks.  
818 *Journal of Computer Science* 10(5), 869–878.
- 819 Liu, H., Hussain, F., Lim, C., Dash, M., 2002. Discretization: An Enabling  
820 Technique. *Data Mining and Knowledge Discovery* 6, 393–423.
- 821 Liu, Y., Zhang, W., Zhang, Z., 2015. A conceptual data model coupling  
822 with physically-based distributed hydrological models based on catchment  
823 discretization schemas. *Journal of Hydrology* 530, 206–215.
- 824 Maldonado, A., Roperó, R. F., Aguilera, P., Rumí, R., Salmerón, A.,  
825 2015. Continuous Bayesian networks for the estimation of species richness.  
826 *Progress in Artificial Intelligence* 4, 49–57.
- 827 Moral, S., Rumí, R., Salmerón, A., 2001. Mixtures of Truncated Exponentials  
828 in hybrid Bayesian networks. In: *ECSQARU’01. Lecture Notes in Artificial  
829 Intelligence*. Vol. 2143. Springer, pp. 156–167.
- 830 Morales, M., Rodríguez, C., Salmerón, A., 2006. Selective naïve Bayes  
831 predictor using mixtures of truncated exponentials. In: *Proceedings of*

- 832 the International Conference on Mathematical and Statistical Modelling  
833 (ICMSM'06).
- 834 Myers, N., Mittenmeier, R. A., Mittenmeier, C. G., da Fonseca, G. A. B.,  
835 Kent, J., 2000. Biodiversity hotspots for conservation priorities. *Nature*  
836 403, 853 – 858.
- 837 Nash, D., Waters, D., Buldu, A., Wu, Y., Lin, Y., Yang, W., Song, Y., Shu,  
838 J., Qin, W., Hannah, M., 2013. Using a conceptual Bayesian network to  
839 investigate environmental management in vegetable production in the Lake  
840 Taihu region of China. *Environmental Modelling & Software* 46, 170–181.
- 841 Newton, A. C., Stewart, G. B., Díaz, A., Golicher, D., Pullin, A. S., 2007.  
842 Bayesian Belief Networks as a tool for evidence-based conservation man-  
843 agement. *Journal for Nature Conservation* 15, 144–160.
- 844 Park, M. H., Stenstrom, M. K., 2008. Classifying environmentally significant  
845 urban land uses with satellite imagery. *Journal of Environmental Manage-*  
846 *ment* 86, 181–192.
- 847 Pollino, C. A., White, A. K., Hart, B. T., 2007. Examination of conflicts and  
848 improved strategies for the management of an endangered eucalypt species  
849 using Bayesian networks. *Ecological Modelling* 201, 37–59.
- 850 Pradhanang, S. M., Briggs, R. D., 2014. Effects of critical source area on  
851 sediment yield and streamflow. *Water and Environmental Journal* 28, 222–  
852 232.
- 853 Roos, T., Wettig, H., Grünwald, P., Myllymäki, P., Tirri, H., 2005. On  
854 discriminative Bayesian network classifiers and logistic regression. *Machine*  
855 *Learning* 59, 267–296.
- 856 Ropero, R. F., Aguilera, P. A., Rumí, R., 2015. Analysis of the socioecological  
857 structure and dynamics of the territory using a hybrid Bayesian network  
858 classifier. *Ecological Modelling* 311, 73–87.
- 859 Rumí, R., 2003. Modelos de redes Bayesianas con variables discretas y con-  
860 tinuas. Ph.D. thesis, Universidad de Almería.
- 861 Rumí, R., Salmerón, A., 2007. Approximate probability propagation with  
862 mixtures of truncated exponentials. *International Journal of Approximate*  
863 *Reasoning* 45, 191–210.



- 864 Rumí, R., Salmerón, A., Moral, S., 2006. Estimating mixtures of truncated  
865 exponentials in hybrid Bayesian networks. *Test* 15, 397–421.
- 866 Schmitz, M., Pineda, F., Castro, H., Aranzabal, I. D., Aguilera, P., 2005.  
867 Cultural landscape and socioeconomic structure. Environmental value and  
868 demand for tourism in a Mediterranean territory. *Consejería de Medio  
869 Ambiente. Junta de Andalucía. Sevilla.*
- 870 Scott, M. W., 2010. *Logistic Regression. From Introductory to Advanced  
871 Concepts and Applications.* SAGE.
- 872 Segurado, P., Araújo, M. B., 2004. An evaluation of methods for modelling  
873 species distribution. *Journal of Biogeography* 31, 1555–1568.
- 874 Shenoy, P. P., West, J. C., 2011. Inference in hybrid Bayesian networks using  
875 mixtures of polynomials. *International Journal of Approximate Reasoning*  
876 52 (5), 641–657.
- 877 Uusitalo, L., 2007. Advantages and challenges of Bayesian networks in envi-  
878 ronmental modelling. *Ecological Modelling* 203, 312–318.
- 879 van der Gaag, L. C., Renooij, S., Feelders, A., de Groote, A., Eijkemans,  
880 M. J. C., Broekmans, F. J., Fauser, B. C. J. M., 2009a. Aligning Bayesian  
881 network classifiers with medical contexts. In: Perner, P. (Ed.), *Proceed-  
882 ings of the Sixth International Conference on Machine Learning and Data  
883 Mining in Pattern Recognition.* Vol. 5632 of *Lecture Notes in Artificial  
884 Intelligence.* Springer-Verlag, Berlin Heidelberg, pp. 787–801.
- 885 van der Gaag, L. C., Renooij, S., Steeneveld, W., Hogeveen, H., 2009b. When  
886 in doubt ... be indecisive. In: Sossai, C., Chemello, G. (Eds.), *Proceedings  
887 of the European Conference on Symbolic and Quantitative Approaches  
888 to Reasoning with Uncertainty.* Vol. 5590 of *Lecture Notes in Artificial  
889 Intelligence.* Springer-Verlag, Berlin Heidelberg, pp. 518–529.
- 890 Voinov, A., Bousquet, F., 2010. Modelling with stakeholders. *Environmental  
891 Modelling & Software* 24, 1268–1281.
- 892 Yang, Y., Webb, G., 2009. Discretization for naive-Bayes learning: Managing  
893 discretization bias and variance. *Machine Learning* 74, 39–74.

- 894 Yang, Y., Webb, G., Wu, X., 2010. Discretization Methods. Springer, Ch.  
895 Data Mining and Knowledge Discovery Handbook, pp. 101–116.
- 896 Zhou, Y., Fenton, N., Neil, M., 2014. Bayesian network approach to multi-  
897 nomial parameter learning using data and expert judgments. International  
898 Journal of Approximate Reasoning 55, 1252–1268.