

A framework for species distribution modelling with improved pseudo-absence generation

Maialen Iturbide^{a,1,*}, Joaquín Bedia^a, Sixto Herrera^{b,c}, Oscar del Hierro^d,
Miriam Pinto^d, Jose Manuel Gutiérrez^a

^a*Meteorology Group, Institute of Physics of Cantabria, Universidad de Cantabria-CSIC.
39005 Santander, Spain*

^b*Meteorology Group, Department of Applied Mathematics and Computer Science,
Universidad de Cantabria, 39005 Santander, Spain.*

^c*Predictia Intelligent Data Solutions, S.L. CDTUC Fase A, Planta 2–203. Avda. los
Castros s/n 39005 Santander, Spain*

^d*NEIKER-Tecnalia, Basque Institute for Agricultural Research and Development. 48160
Derio, Spain.*

Abstract

Species distribution models (SDMs) are an important tool in biogeography and phylogeography studies, that most often require explicit absence information to adequately model the environmental space on which species can potentially inhabit. In the so called *background pseudo-absences* approach, absence locations are simulated in order to obtain a complete sample of the environment. Whilst the commonest approach is random sampling of the entire study region, in its multiple variants, its performance may not be optimal, and the method of generation of pseudo-absences is known to have a significant influence on the results obtained. Here, we compare a suite of classic (random sampling) and novel methods for pseudo-absence

*Corresponding author. Tel.: +34 942 20 20 64

Email address: miturbide@ifca.unican.es (Maialen Iturbide)

¹Edificio Juan Jordá, Avda Los Castros s/n, 39005 Santander, Spain

data generation and propose a generalizable three-step method combining environmental profiling with a new technique for background extent restriction. To this aim, we consider 11 phylogenetic groups of Oak (*Quercus* sp.) described in Europe. We evaluate the influence of different pseudo-absence types on model performance (area under the ROC curve), calibration (reliability diagrams) and the resulting suitability maps, using a cross-validation approach. Regardless of the modelling algorithm used, random-sampling models were outperformed by the methods that incorporate environmental profiling of the background, stressing the importance of the pseudo-absence generation techniques for the development of accurate and reliable SDMs. We also provide an integrated modelling framework implementing the methods tested in a software package for the open source R environment.

Keywords: Ecological niche, *Quercus*, environmental profiling, sampling methods, threshold distance

1 Introduction

Species Distribution Models (SDMs) constitute rules that associate known presence locations of biological entities with the characteristics of their environment to predict its potential distribution in the geographic space (Guisan and Zimmermann, 2000; Elith and et al, 2006). SDM building techniques can be broadly classified into two types: *profile* and *group discrimination* techniques. The first group refers to those modelling approaches that rely solely on known presences to infer the potential distribution of the species, while group discrimination techniques require information of the environmental range where the species do not occur, that is, absence data. Group discrimi-

11 nation techniques have gained popularity in recent years, as they have been
12 reported to yield better results than profile techniques (Engler et al., 2004;
13 Chefaoui and Lobo, 2008; Elith and et al, 2006; Mateo et al., 2010). However,
14 in part due to the great effort involved in true absence sampling, most of the
15 available datasets for predictive modelling (generally natural history collec-
16 tions, see. e.g. Araújo and Williams, 2000) are lacking explicit absence data.
17 Thus, in most cases discrimination techniques are used, requiring the envi-
18 ronmental characterization of the sites of presence in front of a background
19 sample (pseudo-absence data) that characterizes the available environment
20 in the study region.

21 Although the strong influence of the pseudo-absence generation process
22 has been shown in previous studies, comparative analyses addressing the
23 suitability of different methods, some of them quite novel, are scarce in the
24 literature (Zaniewski et al., 2002; Phillips et al., 2009; Lobo et al., 2010),
25 and there is not a consensus on the way in which pseudo-absences should be
26 generated. In fact, several previous studies addressing this issue (e.g. Hengl
27 et al., 2009; Wisz and Guisan, 2009; Stokland et al., 2011; Senay et al., 2013)
28 propose contradictory solutions. As such, the inclusion of reliable pseudo-
29 absences in model evaluation remains an open issue.

30 The most simple and widely applied method of generating pseudo-absences
31 is random selection of the entire study area (e.g., Gastón and García-Viñas,
32 2011; Hanspach et al., 2011; Domisch et al., 2013). A search in the SCOPUS
33 database containing the terms “habitat suitability”, “niche modelling” and
34 “background data”, “pseudo-absence” or “presence-only”, narrowed to the
35 journals of the first quartile and the topic “environmental sciences” for the

36 period 2009–july 2014, yielded a total of 64 articles from which roughly 80%
37 used presence–only datasets. Of them, the 92% used randomly generated
38 pseudo–absences within the study area, either explicitly (38%), or implic-
39 itly (54%) via the MAXENT algorithm (see e.g.: Barbet-Massin et al., 2012;
40 Jiménez-Valverde, 2012, for details), other 28% used profile techniques and
41 a 12% used target group background (note that some of the articles anal-
42 ysed used more than one type of technique, and therefore percentages do
43 not sum up to 100%). Percentages under 10% correspond to the novel ap-
44 proaches analysed in this article. In spite of its wide application, the random
45 sampling method rises the risk of introducing false absences into the model
46 from locations that are suitable for the species, leading to underestimates of
47 its fundamental niche and potential distribution (Anderson and Raza, 2010).
48 This occurs naturally due to biotic interactions and dispersal limitations that
49 do not allow the species to inhabit, and also very often as a result of sampling
50 biases in the data collections. Faced with this problem, it is common practice
51 to set a buffer distance from known presence localities in order to minimize
52 the false negative rate (e.g., Mateo et al., 2010; Bedia et al., 2013). More elab-
53 orated approaches employ a presence–only algorithm as a preliminary step to
54 move pseudo–absences away in the environmental space (see e.g.: Zaniwski
55 et al., 2002; Engler et al., 2004; Barbet-Massin et al., 2012; Liu et al., 2013) or
56 apply a geographically weighted exclusion, which keeps pseudo–absences out
57 from presences using distance maps (Hirzel et al., 2001; Barbet-Massin et al.,
58 2012; Norris et al., 2011; Hengl et al., 2009). These strategies are intended
59 to reduce the background data to those areas where false absences are less
60 likely to occur, while the target group background method has been posited

61 as a solution to remove some of the bias in presence–data collections, using
62 the presence localities of other species as biased background data (Phillips
63 et al., 2009).

64 Another critical matter regarding pseudo–absence data is the extent from
65 which background is sampled. In fact, the available data in the background
66 is usually much larger than the data characterized by presence localities
67 (Anderson and Raza, 2010). A constrained distribution of pseudo–absences
68 around presence locations can lead to misleading models, while unconstrained
69 sampling can artificially inflate test statistics, as well as the weight of less
70 informative response variables (Van der Wal and Shoo, 2009). As a result,
71 the three–step method has been recently proposed as an adequate approach
72 to overcome these limitations, envisaged to define the extent and the envi-
73 ronmental range of the background from which pseudo–absences are sam-
74 pled (Senay et al., 2013, see Sec. 2.4 for details). From an ecological per-
75 spective, the uncertainty associated to the presence of a biological entity is
76 a combined effect of separate factors (biotic, abiotic and movement factors),
77 that in turn depend on the environment of a specific site. In this context, the
78 three–step method pursues the estimation of the fundamental distribution
79 (regions of favorable abiotic factors) by the introduction of pseudo–absences
80 within the niche space corresponding to areas of non-presence (outside the
81 realized niche) and where movement factors are likely favorable (accessible
82 geographic areas) but not so the abiotic factors (Peterson et al., 2011). On
83 the opposite, random sampling would produce predictions closer to a realized
84 distribution, since it only excludes the presence locations for pseudo–absence
85 data generation.

86 The aims of this study are: (i) to analyze the effect of the method used
87 for pseudo-absence data generation on resulting SDMs, and (ii) to provide
88 a modelling framework implementing the state-of-the-art techniques yielding
89 optimal results. In particular, we compare five pseudo-absence data genera-
90 tion methods, ranging from the classical random sampling of the whole region
91 and the target group method, to more sophisticated three-step techniques,
92 combining environmental profiling and spatial restrictions on the sampling
93 domain. We also propose a new criterion for background extent selection
94 based on the theoretical properties of model performance as a function of
95 distance to presence locations. We consider three modelling techniques com-
96 monly used in SDM applications and 11 phylogenetic groups of *Quercus* sp.
97 identified in Europe (*Quercus* sp Europe database, Petit et al., 2002b). In
98 addition, we provide an integrated modelling framework based on the open-
99 source R language (R Core Team, 2014), implementing the methods tested
100 in this study (Supplementary Material).

101 **2. Methods and materials**

102 *2.1. Species Data*

103 The term “species” is a taxonomic designation, and may not necessarily
104 refer to an ecologically homogeneous group of organisms when different eco-
105 types occur within the study area (Oney et al., 2013). Experimental evidence
106 suggests that conventional SDM is not able to properly capture the climatic
107 response of species by treating them as homogeneous units (Beierkuhnlein
108 et al., 2011). With this regard, Hernández et al. (2006) suggested that
109 research in environmental niche modelling should focus on broad distribu-

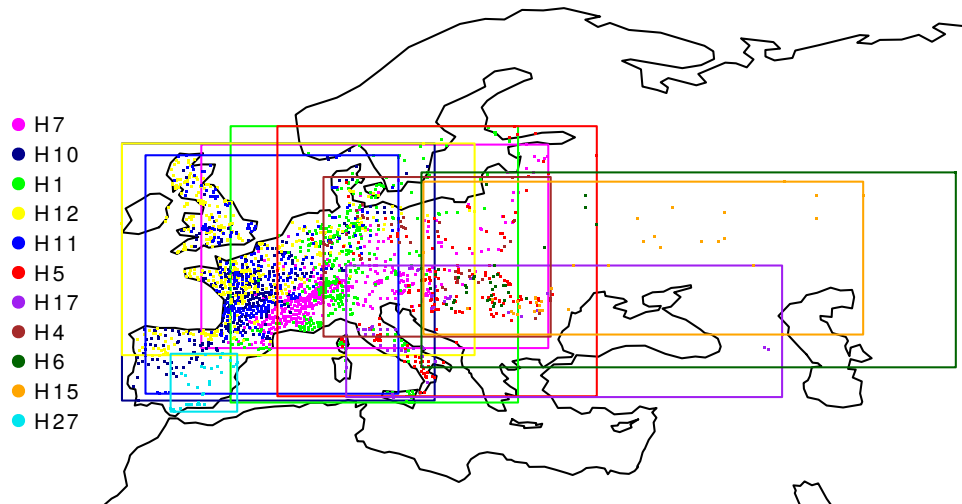


Figure 1: Phylogenetic distribution of *Quercus* sp in Europe. Oak groups in decreasing sample size order are: H7($n=734$), H10($n=651$), H1($n=490$), H12($n=466$), H11($n=283$), H5($n=250$), H17($n=67$), H4($n=53$), H6($n=41$), H15($n=36$) and H27($n=31$).

110 tional subunits based on distinct genetic lineages. For instance, González
 111 et al. (2011) demonstrated that omission error is reduced when “biologi-
 112 cally meaningful” data (in reference to genetically distinct populations of
 113 the same species) are modelled. Hence, in this study we consider genetically
 114 differentiated groups of *Quercus* sp in Europe. Each group corresponds to a dif-
 115 ferent chloroplast haplotype, determined by PCR analysis on more than 2600
 116 populations of Oaks in Europe (see Petit et al., 2002a,b,c). We considered
 117 11 out of the total 42 Oak haplotypes identified, attending to the minimum
 118 population size needed to build the models ($n > 30$) while attending to the
 119 best possible representation of all European *Quercus* lineages (Petit et al.,
 120 2002b, Table 1).

121 The study area was divided in 11 parts (in correspondence to each hap-

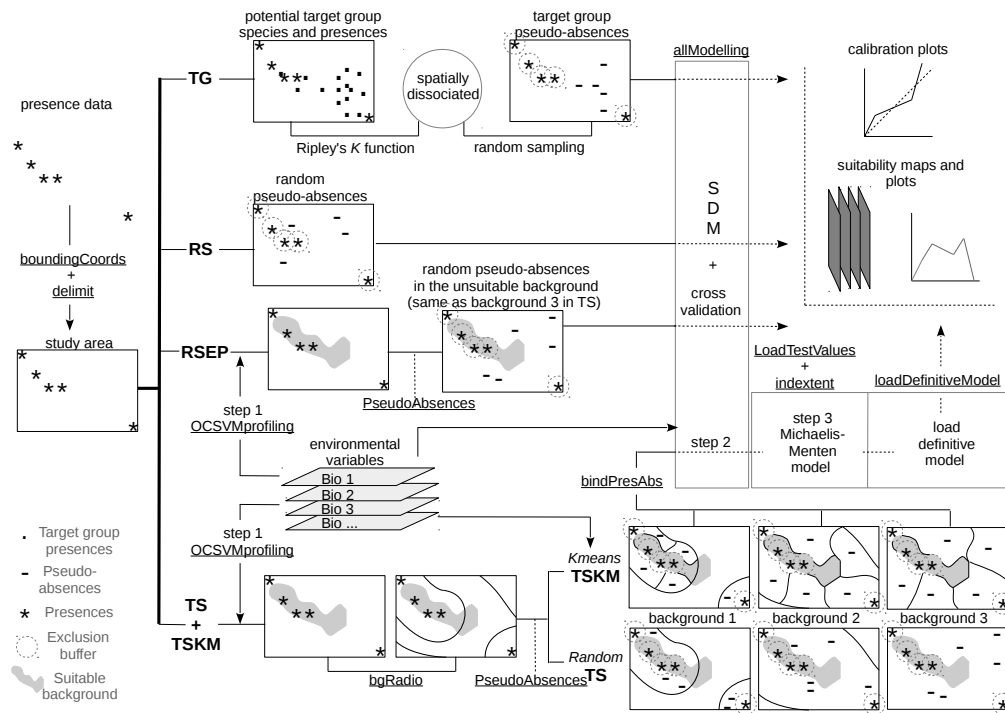


Figure 2: Conceptual diagram of the methodology used in this study. Legend is shown in the bottom left corner. Underlined words refer to the names of the R functions performing each step in the developed framework (see Supplementary Material).

122 lotype distribution) by defining a bounding box around the presence points
 123 (Fig. 1).

124 *2.2. Climate Data*

125 We used the bioclimatic variables of the WorldClim dataset (Hijmans
 126 et al., 2005) at 10 km resolution as explanatory variables to build the SDMs.
 127 The chosen resolution is adequate to the aims of this study, given the ‘false
 128 precision’ provided by the downscaled WorldClim climate surfaces of 1 Km,
 129 as highlighted in previous niche modelling studies (Bedia et al., 2013). After

Table 1: Haplotypes considered ordered by decreasing sample size (n), and the lineages they belong to, according to the *Quercus* sp Europe database (Petit et al., 2002b). Only one lineage (F) out of five was not included in the analyses due to insufficient sample size of all its haplotypes.

Haplotype	Lineage	n
H7	A	734
H10	B	651
H1	C	490
H12	B	466
H11	B	283
H5	A	250
H17	E	67
H4	A	53
H6	A	41
H15	E	36
H27	D	31

130 a pairwise cross-correlation analysis of the bioclimatic variables (following
 131 Bedia et al., 2013), we retained a subset of uncorrelated predictors (bio02,
 132 bio03, bio08, bio13, bio14 and bio15) rescaled in the range [0,1].

133 *2.3. SDM development and assessment*

134 SDMs were built using three different popular techniques, namely max-
 135 imum entropy (MAXENT, Phillips et al., 2006), generalized linear models
 136 (GLMs, Guisan and Zimmermann, 2000) and multivariate adaptive regres-
 137 sion splines (MARS Friedman, 1991). Constrained by data availability, we
 138 resorted to cross-validation techniques (Steyerberg et al., 2010) to replace
 139 truly independent data for model validation, as it is commonplace in ecolog-

140 ical studies (e.g. Manel et al., 1999). In particular, we used a 10-fold cross
141 validation approach, given that it is equally efficient in the error estima-
142 tion as other techniques computationally more demanding like for instance
143 leave-one-out cross validation (Kohavi, 1995).

144 We used the area under the ROC curve (AUC) as the most widely used
145 metric for model performance assessment. The ROC curve describes the pre-
146 dictive ability of the system under the whole range of probability thresholds,
147 thus representing a global measure of model performance, that is quantita-
148 tively assessed by the area it encloses. Thus, high AUC values (closer to 1)
149 indicate good model discrimination, although this is not necessarily coupled
150 to a high numerical accuracy of the predictions (Bedia et al., 2011). With
151 this regard, *calibration plots* (also known as *reliability diagrams*) can be used
152 in order to provide additional information regarding the level of agreement
153 between predicted and observed probabilities of occurrence. This informa-
154 tion is displayed in the form of a plot such that the better the agreement, the
155 closer the line is to the diagonal for the whole range of probability values (see
156 e.g. Bedia et al., 2011; Vaughan and Ormerod, 2005, for a wider explanation
157 in the context of SDM assessment).

158 *2.4. Pseudo–Absence data generation*

159 A larger proportion of pseudo–absences against presences can affect model
160 performance positively or negatively, introducing biases in model inter-comparisons,
161 for which prevalence should be kept constant at an intermediate level (McPher-
162 son et al., 2004; Liu et al., 2005). Thus, for all methods tested we kept
163 the number of pseudo–absences equal to the number of presences in all
164 cases (prevalence = 0.5, Hengl et al., 2009; Mateo et al., 2010; Hanspach

165 et al., 2011; Senay et al., 2013). Additionally, a exclusion buffer of 10 km
166 around the occurrence points was set in order to avoid cells containing both
167 presence and pseudo-absence data (Chefaoui and Lobo, 2008). All steps
168 involved in pseudo-absence generation according to the different methods
169 tested are indicated in the diagram of Fig. 2.

170 *Random selection (RS)*. Pseudo-absences were sampled at random in the
171 whole background, excepting the grid points within the exclusion buffer.

172 *Random selection with environmental profiling (RSEP)*. The RSEP method
173 is aimed at defining the environmental range of the background from which
174 pseudo-absences are sampled. Environmentally unsuitable areas are defined
175 using a presence-only profiling algorithm. To this aim, we run one-class sup-
176 port vector machines (OCSVM, Scholkopf and Smola, 2001) for each Oak
177 group (see e.g. Drake et al., 2006; Bedia et al., 2011, for specific details on the
178 use of support vector machines in SDM studies). OCSVM has been indicated
179 as the most adequate algorithm for this purpose as it can handle high dimen-
180 sional data and complex non-linear relationships between predictors (Senay
181 et al., 2013).

182 *Three-step selection (TS)*. The TS method adds two more steps to the RSEP
183 method to define the environmental range, and also the extent of the back-
184 ground from which pseudo-absences are sampled (Fig. 2). Thus, the first
185 step is the definition of the environmentally unsuitable areas as is done in
186 the RSEP method.

187 In the second step, alternative SDMs are built using random pseudo-
188 absences generated for different spatial extents within the unsuitability back-

189 ground zones defined in the first step. In order to consider all possible extents,
190 we set different maximum *distance thresholds* to each presence location, con-
191 sidering a sequence from 20 km (twice the exclusion buffer) to the length of
192 half diagonal of the bounding box (the maximum possible distance between
193 any pair of points within the area (Fig. 1)), each 10 km (the grid resolution).

194 The third step consists in selecting the optimum background extent and
195 the corresponding fitted model from all possible pseudo-absence configura-
196 tions generated in step 2. Senay et al. (2013) limited the background data
197 using a variable importance change criterion based on principal component
198 analysis to reduce the dimensionality of the environmental space. In our
199 case, we applied a model performance criterion, as variable importance may
200 not always vary significantly for the whole range of distances tested. Thus,
201 a threshold extent is chosen according to the best model performance, while
202 minimizing the distance to presences. With this regard, Van der Wal and
203 Shoo (2009) evaluated the relationship between the geographic extent from
204 which pseudo-absences are taken and model performance, and found that
205 AUC rapidly increased as background size expanded from 10 to 100 km
206 while subsequent expansions resulted in only minor increases in AUC. We
207 found a similar behaviour for all Oak groups, and concluded that the AUC
208 *vs.* distance curve can be optimally fit to an asymptotic Michaelis-Menten
209 type model of the form:

$$v(x) = \frac{Vm \times x}{Km + x}, \quad (1)$$

210 where v and x represent the AUC and the background extent respectively,
211 Vm (Fig. 3) is the asymptotic AUC value achieved by the system and the

212 Michaelis constant Km is the extent at which the AUC is half of Vm . As
213 a result, we propose a generalizable method to find the threshold extent
214 for pseudo-absence sampling near the suitability boundary of the species,
215 without penalizing model performance, which constitutes the major novelty
216 in comparison with previous published methodologies. Thus, AUCs from
217 the multimodel and the different background extents tested are fitted to the
218 curve of equation 1 to extract the theoretical asymptotic AUC value (Vm).
219 Then, the minimum threshold extent x at which $AUC_x > Vm$ is chosen (Fig.
220 3), and the corresponding fitted SDM is retained to produce the suitability
221 maps for the entire study area.

222 *Three-step with k-means selection (TSKM)*. The difference of TSKM with
223 regard to TS is that the pseudo-absences are taken from the spatial subunits
224 defined by a clustering on the background extent in Step 2. Instead of using
225 a random selection on the unsuitable areas after Step 1, a k-means clustering
226 is applied on the environmental and geographical space (k being equal to the
227 number of presence points) and the coordinate values of each cluster centroid
228 are retained, thus obtaining a regular distribution of dissimilar points for
229 the study area which constitutes a representative sample of the unsuitable
230 environment (Senay et al., 2013). Step 3 is then done as in TS method. The
231 resulting background extents for the TS and TSKM methods are listed in
232 Table 2.

233 *Target group selection (TG)*. In order to select a target group for each phylo-
234 genetic Oak group we searched for presence records of species not belonging
235 to the *Fagaceae* family in the database of The Global Biodiversity Infor-
236 mation Facility (GBIF, <http://data.gbif.org>). To ensure a sufficiently

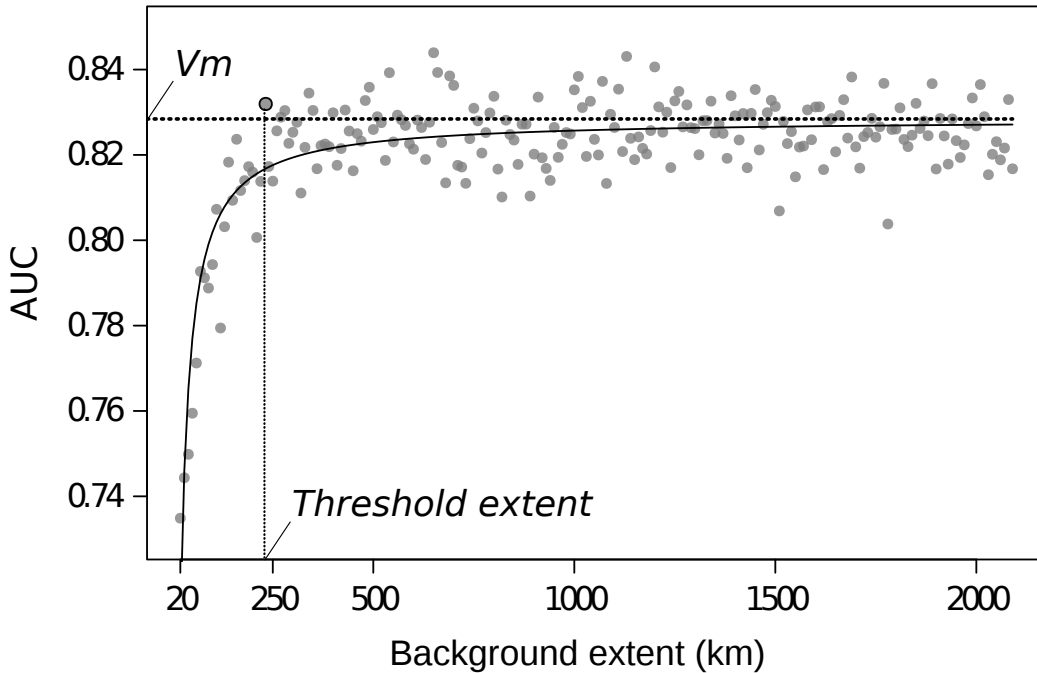


Figure 3: Relation of the AUC to the background extent for group H7. The black curve correspond to the fitted Michaelis-Menten model. V_m represents the maximum AUC achieved by the system. The highlighted point corresponds to the smallest background extent greater than V_m (i.e., the threshold extent). This relationship is similar to that described in Figure 2 in Van der Wal and Shoo (2009). All Oak groups in the study exhibited the same type of curve (see also the examples in the Supplementary Material).

237 high number of presence points, we focused on species with a widespread
 238 distribution in Europe as target group candidates.

239 For each candidate and Oak group, we computed the cross type of the
 240 Ripley's K function (Dixon, 2006) to analyse the spatial behaviour of the
 241 point pattern. From the estimated Cross K -functions, those showing spa-
 242 tial dissociation of the TG candidate with regard to the Oak group were
 243 chosen (see Grantham, 2012, for wider explanation regarding point pattern

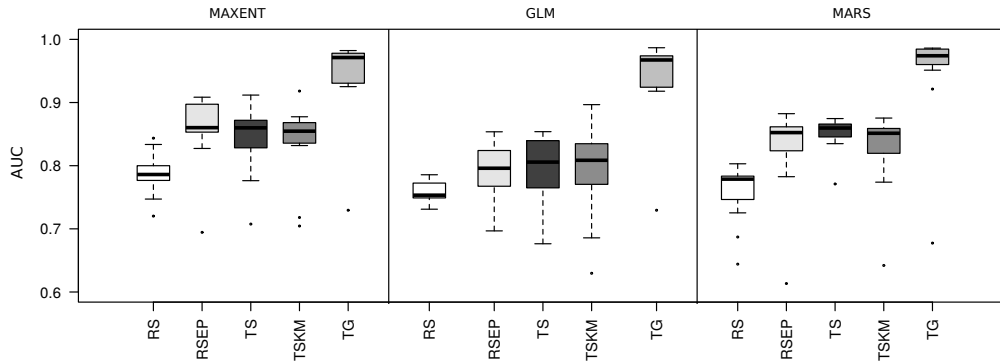


Figure 4: AUC box-plots of the 11 oak groups modelled with the five pseudo-absence generation methods for each modelling technique. Oak groups were modelled with higher accuracy by MAXENT and MARS. The average AUC values improved for all modelling techniques when using a different method from RS.

244 analysis and Rypley’s K function interpretation), resulting in the follow-
 245 ing target groups: *Ulex europaeus* for groups H3 and H11; *Picea glauca* for
 246 groups H1, H2, H4, H5, H6 and H8; *Pinus nigra* for groups H7 and H10;
 247 *Pinus strobus* for group H9. TG locations were then randomly sampled to
 248 match the number of Oak localities in order to obtain balanced datasets for
 249 model training (see Sec. 2.4).

250 3. Results and Discussion

251 3.1. TG method

252 TG attained the highest AUCs for almost all the phylogenetic groups
 253 (Table 3, Fig. 4), but in turn it yielded poorly calibrated models (Fig. 5),
 254 with a strong under-estimation of high probability values. We argue that
 255 these results are due to the spatially clustered distribution of targeted group

Table 2: Threshold distances to presences (kilometres) defining the background extents from which pseudo-absences are sampled. Each data in the column d_{max} correspond to the length of the half diagonal of the bounding box that encloses the study area (Fig. 1), i.e.: the maximum possible distance between a pair of points within the study area.

	d_{TS}	d_{TSKM}	d_{max}
H7	230	290	2090
H10	500	670	2100
H1	580	800	2070
H12	620	620	2130
H11	390	560	1800
H5	190	240	2170
H17	690	830	2360
H4	150	380	1440
H6	1000	1050	2950
H15	360	80	2420
H27	30	70	450

256 presences used as pseudo-absences, leading to spatially autocorrelated back-
 257 ground samples resulting in inflated AUC values (González et al., 2011), and
 258 also to an over-estimated suitability for a large proportion of non-sampled
 259 areas (Figs. 6 and 7), as compared to the other methods. Phillips et al.
 260 (2009) and Mateo et al. (2010) recommended the TG pseudo-absence as the
 261 best method for discrimination, resulting in models with the best predictive
 262 performance. We find the same result, with TG attaining the highest AUC
 263 values, although this comes at the cost of a poor model calibration, and there-
 264 fore we do not recommend this technique if reliable suitability maps are to be
 265 obtained. This stresses the importance of well-distributed presence/absence
 266 data across the environmental and geographical space of the study area in

Table 3: Multimodel mean AUC values, according to the four pseudo-absence generation methods tested, for each of the Oak groups analyzed. Values for TG method are underlined when they are the best of all methods. Values in bold are the maximum AUC values excluding the TG method.

	RS	RSEP	TS	TSKM	TG
H7	0.771	0.834	0.832	0.830	<u>0.981</u>
H10	0.772	0.854	0.851	0.856	<u>0.970</u>
H1	0.764	0.822	0.823	0.820	<u>0.976</u>
H12	0.781	0.839	0.864	0.852	<u>0.971</u>
H11	0.760	0.815	0.842	0.846	<u>0.985</u>
H5	0.786	0.830	0.829	0.828	<u>0.977</u>
H17	0.798	0.847	0.878	0.897	<u>0.935</u>
H4	0.720	0.873	0.835	0.824	<u>0.962</u>
H6	0.802	0.847	0.862	0.859	<u>0.939</u>
H15	0.762	0.668	0.748	0.707	<u>0.941</u>
H27	0.726	0.843	0.741	0.677	0.712

267 order to obtain reliable models (Lobo and Tognelli, 2011).

268 3.2. RSEP, TS and TSKM methods

269 RSEP and three-step methods (TS and TSKM) attained similar results.
 270 As expected, we did not find any significant differences in their AUCs (Fig.
 271 4, Table 3) since both TS and TSKM define a threshold extent based on
 272 the asymptotic AUC value Vm (Fig. 3), close to the expected value of the
 273 maximum distance threshold used by the RSEP method. With this regard,
 274 TS and TSKM methods are preferable than RSEP, since using the theoretical
 275 AUC value given by Vm ensures the selection of a good model, while RSEP
 276 method may result in a sub-optimal model if the last point in the X-axis lies

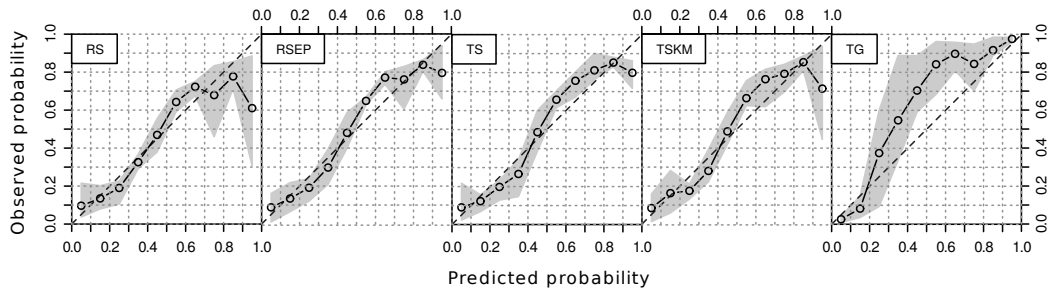


Figure 5: Calibration plots of the multimodel predictions. Points connected by lines are the mean obtained from the different Oak groups and the grey area correspond to the range between maximum and minimum values. Values below the diagonal indicate over-estimated probabilities and values above it under-estimated predictions. The smallest Oak groups H4(n=53), H6(n=41), H15(n=36) and H27(n=31), are excluded in the calibration plots, because their low sample size systematically yields poorly calibrated models that mask observable differences between methods.

277 significantly below the Vm value by chance (Fig. 3).

278 The suitability plots (Fig. 7) show a similar behaviour, clearly differ-
 279 ent from RS and TG. Thus, we conclude that the relevant step that affects
 280 SDM results is the environmental profiling of the background, which con-
 281 stitutes the common characteristic of the RSEP and three-step methods.
 282 As a result, RSEP was equally effective while entailing a more straightfor-
 283 ward implementation. Analogously, since the background extent restriction
 284 does not impair final results, three-step methods are also recommendable
 285 as the effect of non informative pseudo-absences from far regions could be
 286 significant in other case studies, especially when a wider study area is con-
 287 sidered. In this sense, several authors argue that pseudo-absences from far
 288 regions should be avoided (Van der Wal and Shoo, 2009; Anderson and Raza,
 289 2010). Moreover, Jiménez-Valverde (2008) and Lobo et al. (2010) suggested

290 that pseudo-absences should be located near the external boundary of the
291 suitable environment to adequately represent the potential distribution of a
292 species. At this respect, we consider that the three-step method proposed
293 in this study satisfies this requirement while avoids misleading models with
294 reduced AUCs. In addition, TS is generalizable and its implementation is
295 straightforward using the R functions provided (Supplementary Material).
296 Finally, since the TSKM method does not improve SDM results in relation
297 to TS, the introduction of the k-means clustering in Step 2 of TSKM can be
298 skipped in favour of a simple random selection within the background extent.

299 *3.3. RS method vs. RSEP, TS and TSKM methods*

300 The RS method produced well calibrated SDMs, excepting in the zones
301 of higher environmental suitability, where the latter was over-estimated for
302 all Oak groups (Fig. 5). This is due to the fact that many pseudo-absences
303 are distributed around presences inside the potentially suitable environment,
304 resulting in a lower rate of observed presences against absences in the zones
305 predicted as most suitable, and is arguably one major disadvantage of the
306 RS method with regard to methods applying environmental profiling as a
307 previous step (RSEP, TS and TSKM). Furthermore, RS yielded the worst
308 discrimination results, with the lowest AUC values for all algorithms tested
309 (Fig. 4) and for most Oak groups (Table 3).

310 The use of a profiling technique as an intermediate step, characteristic
311 of the three-step methods (TS and TSKM), has been criticized by some au-
312 thors for producing artificially high probabilities of occurrence (Wisn and
313 Guisan, 2009; Stokland et al., 2011) and wider predicted suitability areas.
314 In ecological terms, the variability in the predicted probabilities is related to

315 the ability of the SDMs to represent realized *vs.* potential species distribu-
316 tions, lying spatially wider predicted distributions closer to the fundamental
317 niche of the target species (Chefaoui and Lobo, 2008). However, since the
318 potential distribution of the species is uncertain, we see no reason to pe-
319 nalize the model based on the extent of the area predicted as suitable (see
320 e.g. Jiménez-Valverde, 2012). Furthermore, our results indicate that the pre-
321 dicted potential areas are not significantly shrink/widened with the use of
322 either profiling/RS techniques (they are though in case of TG method, Fig.
323 6). In fact, the most remarkable difference between both is a higher resolu-
324 tion of the profiling-based models as compared to RS for most Oak groups,
325 as depicted by the suitability plots (Fig. 7). This means that ambiguous
326 probabilities (around 0.5) are less likely to occur when RSEP or three-step
327 methods are introduced, in favor of more informative predicted probabilities
328 closer either to 1 or to 0, as opposed to the traditional RS approach. (see
329 e.g. Bedia et al., 2011, for a more detailed explanation of model resolution
330 in the context of SDMs). This is particularly important in order to reduce
331 uncertainties when binary presence/absence maps are required for decision
332 making and/or management plans.

333 Furthermore, the lack of records from suitable regions may simply derive
334 from an inadequate sampling (Anderson, 2003; Hanspach et al., 2011). In
335 fact, presence data is quite often environmentally biased (Bierman et al.,
336 2010) resulting in presence data that does not represent the whole environ-
337 mental range of the realized niche. In these cases, the RS method introduces
338 false absences (within both the realized and fundamental niches) introduc-
339 ing a major source of uncertainty (Lobo et al., 2010) and resulting in over-

340 constrained areas of high suitability (Fig. 7). In this sense, as long as RSEP,
341 TS and TSKM methods sample pseudo-absences within a previously profiled
342 unsuitable area, the risk of introducing false pseudo-absences is minimized,
343 even in the case of relatively biased species collections. On the other hand, in
344 case of error in the initial presence data (e.g. false positives), then profiling
345 techniques may bear the risk of further reinforcing this bias rather than cor-
346 recting it, although this particular situation should be further investigated.

347 *3.4. Sensitivity of model performance to the pseudo-absence generation method*

348 Our results show that the method of pseudo-absence generation strongly
349 conditions output SDMs. Whilst the choice of the SDM algorithm is gen-
350 erally recognized as the principal factor of uncertainty in niche modelling
351 studies (see e.g. Buisson et al., 2010; Fronzek et al., 2011), in this case study
352 we demonstrate that pseudo-absence sampling design is even more impor-
353 tant, leading to a larger variation of model AUC (Fig. 4, Table 3) than
354 the modelling algorithms tested or the initial presence dataset choice, even
355 though MAXENT and MARS performed better than GLMs (Fig. 4), indi-
356 cating that algorithm selection is also an important factor (Phillips et al.,
357 2009; Bedia et al., 2011; Senay et al., 2013). Our results also suggest that
358 MARS performance was more sensitive to the pseudo-absence configuration
359 than MAXENT (Fig. 4), although a more intensive testing beyond the scope
360 of this study would be required to ascertain the sensitivity of different algo-
361 rithms to the pseudo-absence generation scheme.

362 3.5. *Sample size effect on results*

363 As sample sizes are heterogeneous across Oak groups, this allowed us
364 to indirectly evaluate the influence of the sample size in the performance.
365 Caution has to be given to interpreting inflated AUC values due to small
366 number of records (Wisz et al., 2008). For instance, Hanspach et al. (2011)
367 excluded species with less than 50 records to allow reliable modelling. In this
368 study, the calibration analysis shows that group H4 (53 presence records)
369 and smaller groups (Table 1), did not produce reliable models for any of
370 the pseudo-absence generation methods compared (not shown), even though
371 AUC values were generally high (Table 3). In addition, the poor performance
372 of the models for the smallest Oak groups (H15 and H27) is also reflected
373 in the relationship of AUC and background extent, resulting in poor model
374 fits in the TS and TSKM methods (equation 1) and yielding small threshold
375 extents and lower AUCs (Tables 2 and 3).

376 4. **Conclusion**

377 The method for pseudo-absence generation strongly affected output SDM
378 performance regardless of the modelling algorithm chosen and for all the Oak
379 groups tested. The classical random sampling method (RS) yielded the low-
380 est overall performance, while the target group (TG) approach attained high
381 AUC values at the cost of poorly calibrated models, resulting in unreliable
382 suitability maps. Methods that include environmental profiling in a previous
383 step (RSEP, TS and TSKM), clearly outperformed both RS and TG, yield-
384 ing high AUC values and better calibrated predictions, resulting in the most
385 reliable suitability maps with a higher resolution of the predicted probabil-

ities. Thus, we suggest that further investigation on pseudo-absence data
generation should focus in background data profiling. We recommend TS
as the most adequate method, and also RSEP as a computationally simpler
alternative. We also propose the AUC-driven method based on asymptotic
curve fitting as an easily implementable and generalizable approach to ob-
tain a suitable background extent threshold. RSEP, TS and TSKM methods
are implemented in the open source R package `mopa` (*MOdelling Pseudo*
Absences, <https://github.com/miturbide/mopa>), described with worked
examples in the Supplementary Material.

5. Acknowledgments

We are grateful to Rémy Petit and François Ehrenmann for providing
the phylogenetic distribution of *Quercus*. We acknowledge the fruitful dis-
cussions arisen in the the WG1 of the FPS COST Action FP1202 (MaP-FGR,
“Strengthening conservation: a key issue for adaptation of marginal/peripheral
populations of forest trees to climate change in Europe”). We also thank two
anonymous referees for their thoughtful comments that greatly improved the
manuscript. This work was supported by the EC-funded project ADAPTA-
CLIMA II (INTERREG IVB SUDOE Program).

6. References

Anderson, R. P., 2003. Real vs. artefactual absences in species distributions:
tests for *oryzomys albigularis* (rodentia: Muridae) in venezuela. *Journal of*
Biogeography 30, 591–605.

- 408 Anderson, R. P., Raza, A., 2010. The effect of the extent of the study region
409 on GIS models of species geographic distributions and estimates of niche
410 evolution: preliminary tests with montane rodents (genus *Nephelomys*) in
411 Venezuela. *Journal of Biogeography* 37 (7), 1378–1393.
- 412 Araújo, M. B., Williams, P. H., 2000. Selecting areas for species persistence
413 using occurrence data. *Biological Conservation* 96, 331–345.
- 414 Barbet-Massin, M., Jiguet, F., Albert, C. H., Thuiller, W., 2012. Select-
415 ing pseudo-absences for species distribution models: how, where and how
416 many? *Methods in Ecology and Evolution* 3 (2), 327–338.
- 417 Bedia, J., Busqué, J., Gutiérrez, J. M., 2011. Predicting plant species dis-
418 tribution across an alpine rangeland in northern Spain: a comparison of
419 probabilistic methods. *Applied Vegetation Science* 14, 415–432.
- 420 Bedia, J., Herrera, S., Gutiérrez, J. M., 2013. Dangers of using global biocli-
421 matic datasets for ecological niche modeling. limitations for future climate
422 projections. *Global and Planetary Change* 107.
- 423 Beierkuhnlein, C., Thiel, D., Jentsch, A., Willner, E., Kreyling, J., 2011.
424 Ecotypes of European grass species respond differently to warming and
425 extreme drought. *Journal of Ecology* 99, 703–713.
- 426 Bierman, S. M., Butler, A., Marion, G., Kuehn, I., 2010. Bayesian image
427 restoration models for combining expert knowledge on recording activity
428 with species distribution data. *ECOGRAPHY* 33 (3), 451–460.

- 429 Buisson, L., Thuiller, W., Casajus, N., Lek, S., Grenouillet, G., 2010. Un-
430 certainty in ensemble forecasting of species distribution. *Global Change*
431 *Biology* 16, 1145–1157.
- 432 Chefaoui, R. M., Lobo, J. M., 2008. Assessing the effects of pseudo-
433 absences on predictive distribution model performance. *Ecological Mod-*
434 *elling* 210 (4), 478–486.
- 435 Dixon, P. M., 2006. Ripley’s k function. In: *Encyclopedia of Environmetrics*.
436 John Wiley & Sons, Ltd.
- 437 Domisch, S., Kuemmerlen, M., Jähnig, S., Haase, P., 2013. Choice of study
438 area and predictors affect habitat suitability projections, but not the per-
439 formance of species distribution models of stream biota. *Ecological Mod-*
440 *elling* 257, 1–10.
- 441 Drake, J. M., Randin, C., Guisan, A., Jun. 2006. Modelling ecological niches
442 with support vector machines. *Journal of Applied Ecology* 43, 424–432.
- 443 Elith, J., et al, 2006. Novel methods improve prediction of species’ distribu-
444 tions from occurrence data. *Ecography* 29, 129–151.
- 445 Engler, R., Guisan, A., Rechsteiner, L., 2004. An improved approach for
446 predicting the distribution of rare and endangered species from occurrence
447 and pseudo-absence data. *Journal of Applied Ecology* 41 (2), 263–274.
- 448 Friedman, J. H., 1991. Multivariate adaptive regression splines. *The Annals*
449 *of Statistics* 19.

- 450 Fronzek, S., Carter, T., Luoto, M., 2011. Evaluating sources of uncertainty
451 in modelling the impact of probabilistic climate change on sub-arctic palsa
452 mires. *Natural Hazards and Earth System Sciences* 11, 2981–2995.
- 453 Gastón, A., García-Viñas, J., 2011. Modelling species distributions with pe-
454 nalisated logistic regressions: A comparison with maximum entropy models.
455 *Ecological Modelling* 222 (13), 2037–2041.
- 456 González, S., Soto-Centeno, J., Reed, D., 2011. Population distribution mod-
457 els: Species distributions are better modeled using biologically relevant
458 data partitions. *BMC Ecology* 11.
- 459 Grantham, N., 2012. Analyzing multiple independent spatial point processes.
460 *Statistics*.
- 461 Guisan, A., Zimmermann, N. E., 2000. Predictive habitat distribution models
462 in ecology. *Ecological modelling* 135 (2), 147–186.
- 463 Hanspach, J., Kühn, I., Schweiger, O., Pompe, S., Klotz, S., 2011. Geograph-
464 ical patterns in prediction errors of species distribution models. *Global
465 Ecology and Biogeography* 20 (5), 779–788.
- 466 Hengl, T., Sierdsema, H., Radović, A., Dilo, A., 2009. Spatial prediction of
467 species’ distributions from occurrence-only records: combining point pat-
468 tern analysis, ENFA and regression-kriging. *Ecological Modelling* 220 (24),
469 3499–3511.
- 470 Hernández, P. A., Graham, C. H., Master, L. L., Albert, D. L., 2006. The
471 effect of sample size and species characteristics on performance of different
472 species distribution modeling methods. *Ecography* 29 (5), 773–785.

- 473 Hijmans, R. J., Cameron, S. E., Parra, J. L., Jones, P. G., Jarvis, A., 2005.
474 Very high resolution interpolated climate surfaces for global land areas.
475 *International Journal of Climatology* 25, 1965–1978.
- 476 Hirzel, A. H., Helfer, V., Metral, F., 2001. Assessing habitat-suitability mod-
477 els with a virtual species. *Ecological modelling* 145 (2), 111–121.
- 478 Jiménez-Valverde, A., Lobo, JM., Hortal J., 2012. Not as good as they seem:
479 the importance of concepts in species distribution modelling. *Diversity and*
480 *Distributions* 14, 885–890. .
- 481 Jiménez-Valverde, A., 2012. Insights into the area under the receiver oper-
482 ating characteristic curve (AUC) as a discrimination measure in species
483 distribution modelling. *Global Ecology and Biogeography* 21 (4), 498–507.
- 484 Kohavi, R., 1995. A study of cross-validation and bootstrap for accuracy
485 estimation and model selection. In: *Proceedings of the International Joint*
486 *Conference on Artificial Intelligence*. pp. 1137–1143.
- 487 Liu, C., Berry, P. M., Dawson, T. P., Pearson, R. G., 2005. Selecting thresh-
488 olds of occurrence in the prediction of species distributions. *Ecography*
489 28 (3), 385–393.
- 490 Liu, C., White, M., Newell, G., Griffioen, P., 2013. Species distribution mod-
491 elling for conservation planning in victoria, australia. *Ecological Modelling*
492 249, 68–74.
- 493 Lobo, J. M., Jiménez-Valverde, A., Hortal, J., 2010. The uncertain nature of
494 absences and their importance in species distribution modelling. *Ecography*
495 33 (1), 103–114.

- 496 Lobo, J. M., Tognelli, M. F., 2011. Exploring the effects of quantity and
497 location of pseudo-absences and sampling biases on the performance of
498 distribution models with limited point occurrence data. *Journal for Nature*
499 *Conservation* 19 (1), 1–7.
- 500 Manel, S., Dias, J. M., Buckton, S. T., Ormerod, S. J., 1999. Alternative
501 methods for predicting species distribution: an illustration with himalayan
502 river birds. *Journal of Applied Ecology* 36, 734–747.
- 503 Mateo, R. G., Croat, T. B., Felicísimo, A. M., Muñoz, J., 2010. Profile or
504 group discriminative techniques? generating reliable species distribution
505 models using pseudo-absences and target-group absences from natural his-
506 tory collections. *Diversity and Distributions* 16 (1), 84–94.
- 507 McPherson, J. M., Jetz, W., Rogers, D. J., 2004. The effects of species’
508 range sizes on the accuracy of distribution models: ecological phenomenon
509 or statistical artefact? *Journal of Applied Ecology* 41, 811–823.
- 510 Norris, D., Rocha-Mendes, F., Frosini de Barros Ferraz, S., Villani, J.,
511 Galetti, M., 2011. How to not inflate population estimates? spatial den-
512 sity distribution of white-lipped peccaries in a continuous atlantic forest.
513 *Animal Conservation* 14 (5), 492–501.
- 514 Oney, B., Reineking, B., O’Neill, G., Kreyling, J., 2013. Intraspecific varia-
515 tion buffers projected climate change impacts on *Pinus contorta*. *Ecology*
516 *and Evolution* 3 (2), 437–449.
- 517 Peterson, A.T., Soberón, J., Pearson, R.G. and Robert P. Anderson, R.P. and
518 Enrique Martínez-Meyer, E. and Miguel Nakamura, M. *Ecological Niches*

- 519 and Geographic Distributions (MPB-49), 2011. Princeton: Princeton Uni-
520 versity Press.
- 521 Petit, R. J., Brewer, S., Bordács, S., Burg, K., Cheddadi, R., Coart, E.,
522 Cottrell, J., Csaikl, U. M., van Dam, B., Deans, J. D., Espinel, S., Fineschi,
523 S., Finkeldey, R., Glaz, I., Goicoechea, P. G., Jensen, J. S., König, A. O.,
524 Lowe, A. J., Madsen, S. F., Mátyás, G., Munro, R. C., Popescu, F., Slade,
525 D., Tabbener, H., de Vries, S. G. M., Ziegenhagen, B., de Beaulieu, J.-L.,
526 Kremer, A., 2002a. Identification of refugia and post-glacial colonisation
527 routes of european white oaks based on chloroplast DNA and fossil pollen
528 evidence. *Forest Ecology and Management* 156 (1–3), 49–74.
- 529 Petit, R. J., Csaikl, U. M., Bordács, S., Burg, K., Coart, E., Cottrell, J., van
530 Dam, B., Deans, J. D., Dumolin-Lapégue, S., Fineschi, S., Finkeldey, R.,
531 Gillies, A., Glaz, I., Goicoechea, P. G., Jensen, J. S., König, A. O., Lowe,
532 A. J., Madsen, S. F., Mátyás, G., Munro, R. C., Olalde, M., Pemonge, M.-
533 H., Popescu, F., Slade, D., Tabbener, H., Turchini, D., de Vries, S. G. M.,
534 Ziegenhagen, B., Kremer, A., 2002b. Chloroplast DNA variation in euro-
535 pean white oaks: Phylogeography and patterns of diversity based on data
536 from over 2600 populations. *Forest Ecology and Management* 156 (1–3),
537 5–26.
- 538 Petit, R. J., Latouche-Halle, C., Pemonge, M., Kremer, A., 2002c. Chloro-
539 plast DNA variation of oaks in france and the influence of forest frag-
540 mentation on genetic diversity. *Forest Ecology and Management* 156 (1),
541 115–129.
- 542 Phillips, S. J., Anderson, R. P., Schapire, R. E., 2006. Maximum en-

- 543 tropy modeling of species geographic distributions. *Ecological Modelling*
544 190 (3–4), 231–259.
- 545 Phillips, S. J., Dudík, M., Elith, J., Graham, C. H., Lehmann, A., Leathwick,
546 J., Ferrier, S., 2009. Sample selection bias and presence-only distribution
547 models: implications for background and pseudo-absence data. *Ecological*
548 *Applications* 19 (1), 181–197.
- 549 R Core Team, 2014. R: A Language and Environment for Statistical Com-
550 puting. R Foundation for Statistical Computing, Vienna, Austria.
551 URL <http://www.R-project.org/>
- 552 Scholkopf, B., Smola, A. J., 2001. *Learning with Kernels: Support Vector*
553 *Machines, Regularization, Optimization, and Beyond*. MIT Press, Cam-
554 bridge, MA, USA.
- 555 Senay, S. D., Worner, S. P., Ikeda, T., 2013. Novel three-step pseudo-absence
556 selection technique for improved species distribution modelling. *PLoS ONE*
557 8 (8), e71218.
- 558 Steyerberg, E. W., Vickers, A. J., Cook, N. R., Gerds, T., Gonen, M., Obu-
559 chowski, N., Pencina, M. J., Kattan, M. W., 2010. Assessing the perfor-
560 mance of prediction models: a framework for some traditional and novel
561 measures. *Epidemiology (Cambridge, Mass.)* 21 (1), 128–138.
- 562 Stokland, J. N., Halvorsen, R., Stø a, B., 2011. Species distribution mod-
563 elling—Effect of design and sample size of pseudo-absence observations.
564 *Ecological Modelling* 222 (11), 1800–1809.

- 565 Van der Wal, J., Shoo, L. P., 2009. Selecting pseudo-absence data for
566 presence-only distribution modeling: How far should you stray from what
567 you know? *Ecological Modelling* (4), 589–594.
- 568 Vaughan, I. P., Ormerod, S. J., 2005. The continuing challenges of testing
569 species distribution models: Testing distribution models. *Journal of Ap-
570 plied Ecology* 42, 720–730.
- 571 Wisz, M. S., Guisan, A., 2009. Do pseudo-absence selection strategies influ-
572 ence species distribution models and their predictions? an information-
573 theoretic approach based on simulated data. *BMC Ecology* 9 (1), 8.
- 574 Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan,
575 A., Group, N. P. S. D. W., 2008. Effects of sample size on the performance
576 of species distribution models. *Diversity and Distributions* 14 (5), 763–773.
- 577 Zaniwski, A. E., Lehmann, A., Overton, J. M., 2002. Predicting species
578 spatial distributions using presence-only data: a case study of native new
579 zealand ferns. *Ecological Modelling* 157 (2), 261–280.

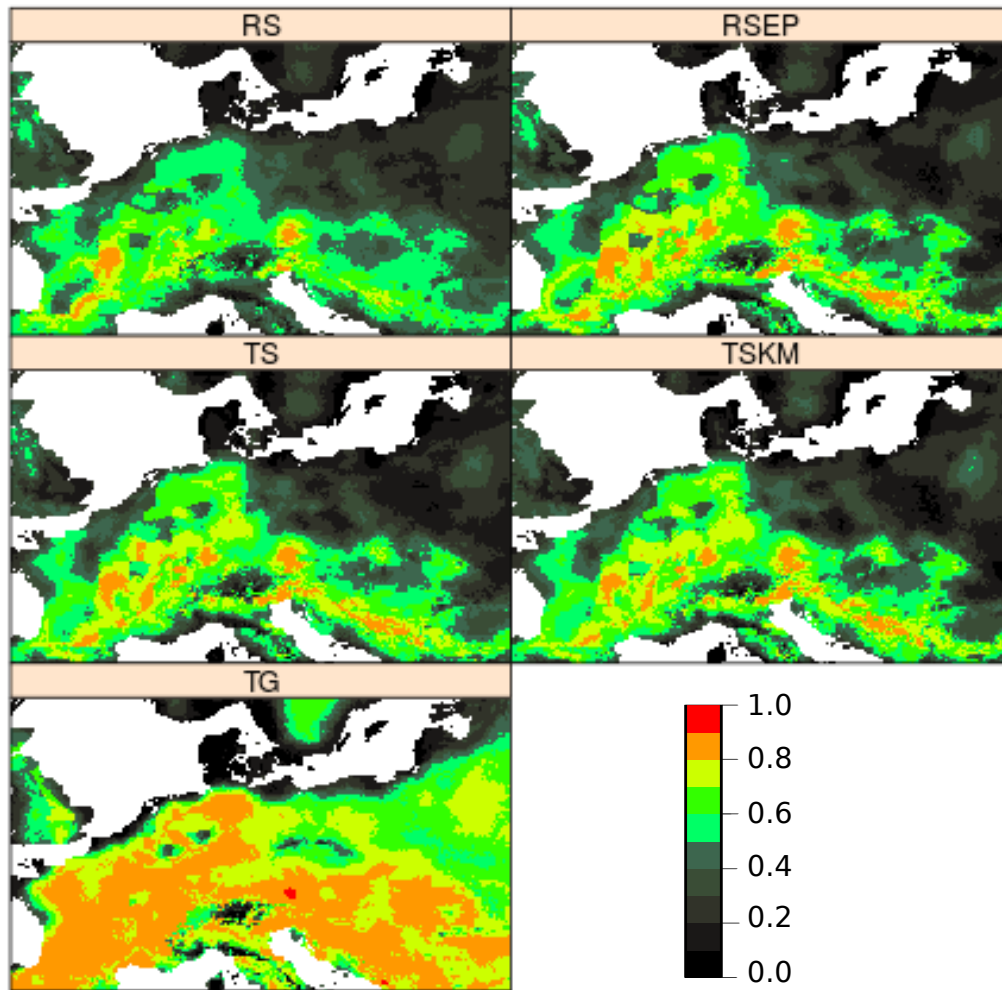


Figure 6: Multimodel suitability maps according to the five pseudo-absence generation methods tested for Oak group H7. Maps for the rest Oak groups show the same pattern on the prediction change between methods as is shown in Figure 7. Suitability is here expressed as a probability of occurrence given the environmental conditions, in the range $[0,1]$.

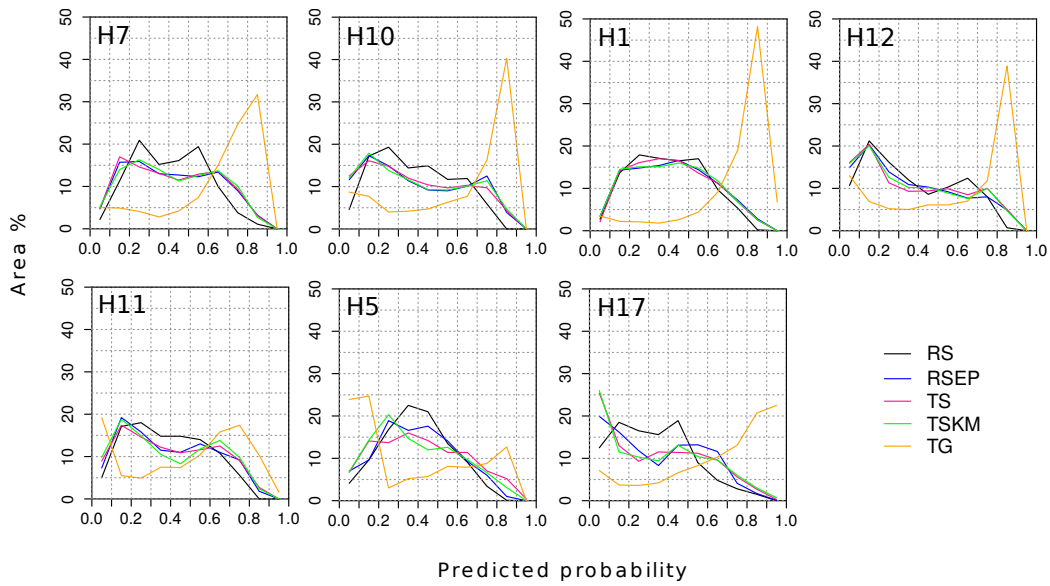


Figure 7: Suitability plots. Percentage of area predicted into each interval of probability of occurrence for the Oak groups producing well calibrated models (see Figure 5). These graphics give quantitative information on the suitability maps for a better interpretation of the results obtained. The first plot (H7) correspond to the suitability maps shown in Figure 6. Compared to RS, the RSEP, TS and TSKM methods produce incremented areas of high and low suitability and reduced mid suitable areas. The TG method predicts large areas of high suitability.

Supplementary material for online publication only

[Click here to download Supplementary material for online publication only: supplementary_material_mopa.pdf](#)

LaTeX Source Files

[Click here to download LaTeX Source Files: Iturbide_manuscript_source_files.zip](#)