

1 **How to best threshold and validate stacked species assemblages? Community**  
2 **optimisation might hold the answer**

3 Daniel Scherrer<sup>a</sup>, Manuela D’Amen<sup>a</sup>, Rui F. Fernandes<sup>a</sup>, Rubén G. Mateo<sup>a,b</sup> and Antoine  
4 Guisan<sup>a,c</sup>

5 <sup>a</sup> Department of Ecology and Evolution, University of Lausanne, Biophore, CH-1015  
6 Lausanne, Switzerland

7 <sup>b</sup> ETSI de Montes, Forestal y del Medio Natural, Universidad Politécnica de Madrid, Ciudad  
8 Universitaria s/n, 28040, Madrid, Spain.

9 <sup>c</sup> Institute of Earth Surface Dynamics, University of Lausanne, Géopolis, CH-1015 Lausanne,  
10 Switzerland

11

12 Corresponding author:

13 Daniel Scherrer

14 Department of Ecology and Evolution

15 University of Lausanne

16 CH-1015 Lausanne, Switzerland

17

18 Running title: Thresholding stacked assemblages

19 **ABSTRACT**

- 20 1. The popularity of species distribution models (SDMs) and the associated stacked  
21 species distribution models (S-SDMs), as tools for community ecologists, largely  
22 increased in recent years. However, while some consensus was reached about the best  
23 methods to threshold and evaluate individual SDMs, little agreement exists on how to  
24 best assemble individual SDMs into communities, i.e. how to build and assess S-SDM  
25 predictions.
- 26 2. Here, we used published data of insects and plants collected within the same study  
27 region to test (1) if the most established thresholding methods to optimize single  
28 species prediction are also the best choice for predicting species assemblage  
29 composition, or if community-based thresholding can be a better alternative, and (2)  
30 whether the optimal thresholding method depends on taxa, prevalence distribution  
31 and/or species richness. Based on a comparison of different evaluation approaches we  
32 provide guidelines for a robust community cross-validation framework, to use if  
33 spatial or temporal independent data are unavailable.
- 34 3. Our results showed that the selection of the “optimal” assembly strategy mostly  
35 depends on the evaluation approach rather than taxa, prevalence distribution, regional  
36 species pool or species richness. If evaluated with independent data or reliable cross-  
37 validation, community-based thresholding seems superior compared to single species  
38 optimisation. However, many published studies did not evaluate community  
39 projections with independent data, often leading to overoptimistic community  
40 evaluation metrics based on single species optimisation.
- 41 4. The fact that most of the reviewed S-SDM studies reported over-fitted community  
42 evaluation metrics highlights the importance of developing clear evaluation guidelines

43 for community models. Here, we move a first step in this direction, providing a  
44 framework for cross-validation at the community level.

## 45 INTRODUCTION

46 Past and future environmental changes may not only lead to shifts in species distributions  
47 (e.g., Parmesan & Yohe 2003; Thuiller *et al.* 2005; Dullinger *et al.* 2012), but also to changes  
48 in species assemblages and interactions (e.g., Van der Putten, Macel & Visser 2010; Nogues-  
49 Bravo & Rahbek 2011; Blois *et al.* 2013; Alexander *et al.* 2016). Information about  
50 communities, here defined as a taxonomic assemblage of distinct populations of species that  
51 co-occur in a given space at a given time (Begon, Harper & Townsend 1996), is therefore  
52 essential to make informed decisions for conservation prioritisation (D'Amen *et al.* 2011;  
53 Guisan *et al.* 2013; Mateo *et al.* 2013) and to create biodiversity indices (e.g., Essential  
54 Biodiversity Variables; Pereira *et al.* 2013) for policy decisions (Fleishman, Noss & Noon  
55 2006; Granger *et al.* 2015).

56 Different approaches to model communities are available, using either correlative (e.g.,  
57 Ferrier & Guisan 2006; Guisan & Rahbek 2011) or mechanistic techniques (e.g., Kearney &  
58 Porter 2009; Mokany & Ferrier 2011), with some predicting only macro-ecological properties  
59 such as species richness (e.g., Currie *et al.* 2004; Gotelli *et al.* 2009; Dubuis *et al.* 2011) and  
60 others also predicting community composition (see D'Amen *et al.* 2017b for a review). In this  
61 study, we focused on correlative approaches based on individual species distribution models  
62 (SDMs), as they are the most common technique applied to conservation strategies (Guisan *et*  
63 *al.* 2013), and to predict future patterns of biodiversity in the face of global change (Nogues-  
64 Bravo & Rahbek 2011; D'Amen *et al.* 2017b). Niche-based SDMs quantify the relationship  
65 between available species occurrences and different environmental factors to analyse and  
66 predict distributional patterns (Guisan & Thuiller 2005; Elith & Leathwick 2009; Guisan,  
67 Thuiller & Zimmermann 2017). By additionally stacking individual SDMs (S-SDMs), one  
68 can produce spatiotemporal projections of species richness and composition (Ferrier & Guisan  
69 2006; Guisan & Rahbek 2011).

70 While there is a vast and now long-standing literature on advances and limitations of single  
71 species predictions (e.g., Guisan & Thuiller 2005; Guisan *et al.* 2006; Maggini *et al.* 2006;  
72 Elith & Leathwick 2009; Meier *et al.* 2010; Zimmermann *et al.* 2010; Merow *et al.* 2014),  
73 studies exploring how to improve community predictions based on aggregated information  
74 from individual SDMs emerged more recently (e.g., Mateo *et al.* 2012; Benito, Cayuela &  
75 Albuquerque 2013; Cord *et al.* 2014; Mod *et al.* 2015; but see Ferrier *et al.* 2002). A  
76 fundamental difference among the proposed solutions is whether to maintain the information  
77 on species composition in the final predictions. For instance, the simple sum of probabilities  
78 of individual SDM predictions usually gives better estimates of species richness, but the  
79 information on species identity is lost (Dubuis *et al.* 2011; Calabrese *et al.* 2014). Therefore,  
80 predictions of community composition have mainly been achieved so far by thresholding the  
81 individual continuous SDM predictions (e.g., probability or suitability index) to obtain binary  
82 maps (Liu, White & Newell 2013) and then stacking the latter at the assemblage level (e.g.,  
83 Pottier *et al.* 2013; D'Amen *et al.* 2015; D'Amen, Pradervand & Guisan 2015).

84 There are several examples in the literature of optimizing thresholding methods for single  
85 species predictions (e.g., Liu *et al.* 2005; Jimenez-Valverde & Lobo 2007; Freeman & Moisen  
86 2008; Liu, White & Newell 2013). These led to a mounting consensus about the most  
87 appropriate methods, with the majority of SDM studies published nowadays using either an  
88 approach maximising the true skills statistics (Max.TSS) or based on the curve in a receiver  
89 operating characteristic plot (Opt.ROC, related to AUC) (see Guisan, Thuiller &  
90 Zimmermann 2017; Table S1). However, the threshold selection can strongly influence the  
91 reliability of the predicted richness and composition of S-SDMs assemblages (Pineda & Lobo  
92 2009; Benito, Cayuela & Albuquerque 2013). It is thus relevant to explore which thresholding  
93 approach provides the best performance in assemblage estimates, and if alternatives exist that  
94 can improve the assemblage prediction from individual SDMs.

95 Studies focussing on S-SDMs tend to over-predict species richness when based on  
96 (thresholded) binary predictions (e.g., Pineda & Lobo 2009; Dubuis *et al.* 2011; Mateo *et al.*  
97 2012; Pottier *et al.* 2013; Pouteau *et al.* 2015), with some exceptions (e.g., D'Amen,  
98 Pradervand & Guisan 2015; Distler *et al.* 2015). Different factors have been proposed to  
99 explain this over-prediction: (1) a statistical bias in thresholding site-level occurrence  
100 probabilities for each species (Calabrese *et al.* 2014); (2) the implicit assumption of  
101 unsaturated communities not assuming an ecological limit for species numbers in assemblages  
102 (environmental carrying capacity; Guisan & Rahbek 2011); (3) the lack of considering  
103 different constraints on community composition (i.e., ecological, evolutionary, historical, or  
104 biological biodiversity drivers; see Mateo, Mokany & Guisan 2017).

105 The commonly used approach to get binary maps from continuous SDM predictions is to use  
106 a species-specific threshold, i.e. each species has a single threshold across all sites ("species  
107 threshold", Calabrese *et al.* 2014). Recently, another community-based approach, called  
108 probability ranking rule (PRR), was proposed to predict assemblage composition from  
109 individual SDMs (D'Amen *et al.* 2015). This method does not require a species-specific  
110 threshold, therefore preventing over-prediction, but site-by-site ecological constraints (e.g.,  
111 macro-ecological models) are applied to assemblages to predict species richness ("site-  
112 threshold").

113 Surprisingly, studies aiming to test and improve S-SDM have used very different approaches  
114 to evaluate the predicted assemblages (Cord *et al.* 2014; Hespanhol *et al.* 2015; Pouteau *et al.*  
115 2015; Thuiller *et al.* 2015; Zurell *et al.* 2016) and this evaluation aspect of the community  
116 modelling procedure has not yet received all the attention it deserves. In most studies,  
117 assemblage predictions are not adequately evaluated because the data used for the evaluation  
118 were already used for individual model fitting, not allowing anymore a correct cross-  
119 validation at the community level. Ideally, the best evaluation method should use spatial or

120 temporal independent data (Elith *et al.* 2006; Guisan, Thuiller & Zimmermann 2017), but if  
121 not available, an appropriate cross-validation approach should at least be set up.

122 Here, we used published high-resolution data of insects (butterflies and grasshoppers) and  
123 plants (forests and grasslands sites), collected within the same study region to (1) test if the  
124 most established thresholding methods for optimal single species prediction (i.e., Max.TSS  
125 and Opt.ROC) are also the best choice for species assemblages, (2) investigate if the optimal  
126 thresholding method depends on taxa, prevalence distribution (Allouche, Tsoar & Kadmon  
127 2006), and/or species richness and (3) provide guidelines for a correct community cross-  
128 validation framework, to be used if spatially- or temporally- independent data are unavailable.

## 129 MATERIALS AND METHODS

### 130 Community data and environmental variables

#### 131 *Study area*

132 The data on all taxa were collected within the same study area located in the western Swiss  
133 Alps of the canton Vaud (Fig. 1; 46°10' to 46°30' N; 6°50' to 7°10' E), covering an area of ca.  
134 700 km<sup>2</sup>, with elevation ranging from 375 to 3210 m a.s.l. and forested areas up to 1900 m a.s.l.  
135 For centuries, agriculture (farming and pasturing) has maintained grasslands among forests and  
136 altered the position of the treeline. The highly variable topography and diverse land use of the  
137 study area, in combination with our high-resolution environmental data (25 x 25 m cell size),  
138 provide a huge range of complex species-environment relationships to test our modelling  
139 framework.

#### 140 *Plant data*

141 The forest data were part of a forest inventory of the canton Vaud conducted between 1988 and  
142 2002 (mostly 1990 to 1994) and consisted of 3076 sites. The forest sites were distributed on a  
143 400 m grid all across the forested area of the canton and had a circular area of 314 m<sup>2</sup> (Fig. 1;  
144 for details see Hartmann, Fouvy & Horisberger 2009). In total, 703 plant species were recorded,  
145 but only 312 (44%) had enough occurrence data (> 20 occurrences) across the dataset for  
146 modelling purposes (see Table 1 for more detailed statistics on the datasets).

147 The grassland dataset was collected between 2002 and 2009 following an equal random-  
148 stratified sampling of non-forested areas in the study area. In total, 911 vegetation sites of 4 m<sup>2</sup>  
149 were sampled (Fig. 1; for more information see Dubuis *et al.* 2011). A total of 905 plant species  
150 were recorded but only the 212 most frequent (>20 occurrences) were selected for modelling  
151 (Table 1).

152 To predict the distribution of the plant species we used five environmental variables: growing



153 degree-day (above 0 °C), moisture index over the growing season (difference between  
154 precipitation and potential evapotranspiration), the sum of potential solar radiation over the  
155 year, slope (in degrees), and topographic position (unit-less, indicating the ridges and valleys).  
156 All these variables were at a 25 m resolution and have been shown to be useful predictors for  
157 plant species in mountain environments (see Dubuis *et al.* 2011; D'Amen *et al.* 2015; Scherrer  
158 *et al.* 2017 for details on predictors).

### 159 *Insect data*

160 Data on butterflies and grasshoppers were respectively collected in 192 and 202 squares of 50  
161 m x 50 m across all the elevational range of the study area (Fig. 1; see Pellissier *et al.* 2012;  
162 Pradervand *et al.* 2013, for more information). In total, 131 butterfly and 41 grasshopper  
163 species were observed, but due to model limitations only the most common 67 butterfly and  
164 20 grasshopper species ( $\geq 20$  occurrences) were considered for modelling (Table 1).

165 For our SDMs we used the same predictors as D'Amen, Pradervand and Guisan (2015): four  
166 bioclimatic variables (solar radiation, summer temperature, annual degree-days and annual  
167 average number of frost days during the growing season), an index of vegetation productivity,  
168 i.e. normalized difference vegetation index (as proxies for trophic resources), and the distance  
169 to forest. These variables were selected as they are not highly correlated ( $< 0.7$ ; Dormann *et al.*  
170 2013) and considered ecologically important for insects (e.g., Turner, Gatehouse & Corey  
171 1987; Hawkins & Porter 2003).

### 172 **The modelling framework**

173 Our modelling framework used three different S-SDM based community modelling pathways  
174 (“single species cross-validation”, “independent data” and “community cross-validation”)  
175 representing the most commonly reported practices in the literature (see Fig. 2 and  
176 “Evaluating community predictions” section).

177 *Single species modelling, thresholding and evaluation*

178 Individual species models were run by generalised linear models (GLM; McCullagh & Nelder  
179 1989), generalised additive models (GAM; Hastie & Tibshirani 1990), random forest (RF;  
180 Breiman 2001) and boosted regression trees (BRT; Elith, Leathwick & Hastie 2008). Models  
181 for species with more than 50 occurrences were fitted by simple SDMs using all five selected  
182 predictors, followed by a weighted (AUC) ensemble forecast (Marmion *et al.* 2009). Species  
183 having only between 20 and 50 occurrence records were fitted by an ensemble bivariate  
184 approach optimised for rare or under-sampled species (Lomba *et al.* 2010; Breiner *et al.*  
185 2015): individual models were calibrated on bivariate combinations of the selected predictors  
186 with all four modelling techniques, followed by a consensus forecast from all the resulting  
187 “small models” weighted by their AUC scores. We used a repeated split-sample procedure  
188 (N=25) for model evaluation, followed by a weighted (AUC) ensemble forecast (across  
189 techniques and split-sample runs).

190 The projected probability outputs of the ensemble models were binarised using two  
191 thresholding schemes: (1) *species-specific-thresholds* (a single threshold calculated for each  
192 species) and (2) *site-specific-thresholds* (differing for each site on the basis of additional  
193 community information, i.e. species richness predictions). We selected seven different  
194 species-specific-thresholding techniques, which can be classified in four major groups: single-  
195 index based, sensitivity and specificity combined, model-building data-only-based, and  
196 predicted probability-based (see Table S1; Liu *et al.* 2005; Nenzen & Araujo 2011 for details  
197 on classification). As the thresholding techniques showed minimal within-group variance (see  
198 Figure S1 and S2), we decided to only present the results for one thresholding technique per  
199 group in the main manuscript. The chosen techniques were: Cohen’s Kappa maximization  
200 approach (*Max.Kappa*; single-index based), TSS maximization approach (*Max.TSS*,  
201 sensitivity and specificity combined), observed prevalence (*Obs.Preval*; model-building data-

202 only-based approach), and average probability approach (*AvgProb*; predicted probability-  
203 based approach; for details on techniques see Table S1). In addition, we applied two site-  
204 thresholds (community-based approaches) using species richness (SR) predictions in  
205 combination with a probability ranking rule (PRR). These methods selected a number of  
206 species equal to the predicted SR on the basis of decreasing probabilities of presence  
207 calculated by the SDMs (D'Amen *et al.* 2015; D'Amen, Pradervand & Guisan 2015).  
208 Therefore, the species with the highest probabilities in a site are selected (considered present)  
209 in decreasing order until the SR predicted for the site is reached. The SR predictions were  
210 derived by either summing the per site probabilities of individual SDMs, obtaining a  
211 prediction of richness for each site (pS-SDM; Dubuis *et al.* 2011) or by a macro-ecological  
212 model (MEM; see D'Amen, Pradervand & Guisan 2015 for details), directly modelling the  
213 richness of the sites. As results from the two site-thresholds were concordant, we only show  
214 here the former (*pS-SDM+PRR*).

215 To evaluate the threshold independent performance of our individual species models, the area  
216 under the curve of a Receiver-Operating Characteristic (ROC) plot (AUC; Fielding & Bell  
217 1997) was calculated based on a repeated split sampling cross-validation (Thuiller, Georges &  
218 Engler 2013). Additionally, based on our independent/cross-validation data we calculated five  
219 threshold dependent metrics for each thresholding technique: the overall accuracy (PCC; i.e.  
220 proportion of correctly classified presence and absences; Fielding & Bell 1997), sensitivity  
221 (proportion of correctly predicted presences), specificity (proportion of correctly predicted  
222 absences), the true skill statistic (i.e. [(sensitivity + specificity) - 1]; TSS; Allouche, Tsoar &  
223 Kadmon 2006) and Cohen's Kappa (Kappa; i.e., overall accuracy but corrected for chance  
224 performance; Cohen 1968).

225 *Evaluating community predictions*

226 All the community predictions were built by stacking binary SDMs of individual species (S-  
227 SDMs; Dubuis *et al.* 2011; Guisan & Rahbek 2011). The three modelling pathways (Fig. 2)  
228 were identical regarding the modelling procedure for single species, thresholding and  
229 community assemblage and only varied in the selection of the data for community calibration  
230 and evaluation.

231 - The “single species cross-validation” (SSCV) approach (Fig. 2) has not fully  
232 “unused/independent” data for community evaluation (i.e. sites not used for the  
233 calibration of any single species). Here, in the process of the cross-validation of all  
234 individual SDMs (i.e. across all species), different sites are selected at each resampling  
235 iteration and for each species, so that all sites are most likely used in at least one split-  
236 sampling run and their information incorporated in the final ensemble model. This  
237 approach cannot thus be considered based on fully independent data. The SSCV  
238 approach has been to date the most common way to model and evaluate communities  
239 predictions based on S-SDMs (Fig. 2; e.g., Dubuis *et al.* 2011; Calabrese *et al.* 2014;  
240 D'Amen, Pradervand & Guisan 2015; Distler *et al.* 2015). As no independent data is  
241 set aside for community evaluation, this approach usually gets evaluated with all the  
242 sites used for calibration. However, to avoid bias in the results due to different  
243 numbers of evaluation sites, we evaluated the SSCV approach only on 30% of the  
244 available sites (identical to the ID and CCV approach below).

245 - The (spatial or temporal) “independent data” (ID) approach (Fig. 2) starts with two  
246 completely independent datasets. One is used for the calibration of the SDMs (i.e.  
247 70% of the sites) and the other set is used (only) to evaluate the performance of the  
248 community predictions (i.e. 30% of the sites; Fig 2; e.g., Benito, Cayuela &  
249 Albuquerque 2013; Pottier *et al.* 2013; Cord *et al.* 2014; D'Amen *et al.* 2015; Zurell *et*  
250 *al.* 2016).

251 - The “community cross-validation” (CCV) approach (Fig. 2) uses a repeated split  
252 sampling of sites (100 repetitions) dividing the available sites into calibration (70%)  
253 and evaluation sets (30%) to perform all the modelling procedure from the single  
254 species prediction to the community assembly (Fig. 2). In contrast to the previous ID  
255 pathway (above), which only uses one (spatial or temporal) fixed independent  
256 evaluation dataset, in the CCV approach all SDMs are fitted at each split-sample  
257 iteration using the same training and test sets for all species, thus minimizing the risk  
258 of bias in the evaluation data (i.e. if the training and test sets differ across species, as  
259 in the ID approach). This repeated cross-validation also allows the  
260 estimation/simulation of confidence intervals for community predictions instead of  
261 just a single value per community. To our knowledge, no study used this community  
262 cross-validation method so far.

263 To compare the community model performance among thresholding techniques and  
264 modelling pathways, we calculated eight different community agreement metrics: 1) the  
265 deviation of the predicted from the observed species richness (SR.deviation), 2) the  
266 proportion of species correctly predicted as present (community sensitivity), 3) the proportion  
267 of species correctly predicted as absent (community specificity), 4) community accuracy  
268 (PCC; i.e. the percent correctly classified species, present or absent), 5) the community TSS  
269 (here measured for a site across all species, rather than for a species across all sites as in  
270 single SDM evaluation; Pottier *et al.* 2013) , 6) the community kappa (same as for TSS, for a  
271 site across species; Pottier *et al.* 2013), and 7) the Sørensen similarity (Sørensen 1948).

#### 272 *Correlation of single species and community evaluation metrics*

273 For each combination of dataset, modelling pathway and thresholding method ( $4 \times 3 \times 9 =$   
274 108) we calculated the average evaluation metric for all five single species metrics and all  
275 seven community metrics. We then calculated the Spearman correlation of all possible

276 combinations of our five single species and seven community evaluation metrics. The  
277 resulting correlation matrix tells us if methods (modelling pathways or thresholding methods)  
278 that yield the highest scores in a certain single species metric also yield the highest score in  
279 the corresponding community evaluation metric.

## 280 **RESULTS**

### 281 **Performance of individual SDMs**

282 As expected the evaluation scores of the individual SDMs were similar to earlier studies  
283 published with the same data (D'Amen *et al.* 2015; D'Amen, Pradervand & Guisan 2015;  
284 Scherrer *et al.* 2017) and their performance was not affected by the chosen community  
285 evaluation approach (Table 1, Table S3). Despite their differences in site SR, prevalence  
286 distribution and species pool the average performance of individual SDMs was similar across  
287 all taxa (Table 1, Table S3). Additionally, the often reported effect of species prevalence on  
288 model performance was only marginal in our study, with rare and common species having  
289 similar average model performance within a given taxonomic group (Fig. S3).

### 290 **Correlation of single species and community evaluation metrics**

291 The correlation between the single species and corresponding community metrics was highest  
292 ( $\text{cor} > 0.93$ ; Table 2) for some combinations of metrics based on partial information from the  
293 contingency table comparing predictions to observations (i.e. PCC, specificity and sensitivity)  
294 and considerably lower for the metrics accounting for all dimensions of the contingency table,  
295 such as TSS and Cohen's Kappa ( $\text{cor} = 0.73$ ; Table 2). Correlations between non-  
296 corresponding single species and community metrics (i.e. Sørensen and SR deviation) tended  
297 to be even lower, with the exception of Kappa versus Sørensen (Table 2).

### 298 **Species richness and compositional similarity**

299 The deviation in species richness between observed and predicted communities was strongly  
300 dependent on the chosen thresholding method (Fig. 3). The thresholding approach that uses  
301 the average predicted probability (*AvgProb*) showed the highest amount of over-prediction  
302 followed by the combined sensitivity and specificity approach (*Max.TSS*). The other three  
303 thresholding methods (*Obs.Preval*, *Max.Kappa* and *pS-SDM+PRR*) performed very similar  
304 and showed overall no tendency to over-predict species richness. There were no significant  
305 differences between the three modelling pathways for any of the studied taxa (Fig. 3). The  
306 absolute number of over-predicted species was strongly related to the average number of  
307 species per plot (SR) and therefore differed among the taxa (Fig. 3). However, when corrected  
308 for the differences in SR the over-prediction did not significantly vary anymore across taxa.

309 The compositional similarity (Sørensen similarity index) varied significantly both among  
310 thresholding techniques and modelling pathways (Fig. 4). The compositional similarity was  
311 expectedly always much higher with the “single species cross-validation” (SSCV) pathway  
312 compared to the “independent data” (ID) or the “community cross-validation” (CCV)  
313 pathways, which both performed similarly. There was also a strong interaction between  
314 modelling pathway and thresholding technique. Using the SSCV pathway, thresholding by  
315 *Obs.Preval* and by *Max.Kappa* performed better (Fig. 4). However, if independent sites were  
316 available for the community evaluations (ID and CCV pathways), the community based  
317 approaches (*pS-SDM+PRR*) performed better than the *Obs.Preval* and *Max.Kappa* thresholds  
318 (Fig. 4). The similarity between predicted and observed communities was higher in the two  
319 insect datasets than in the two plant datasets (Fig. 4), which is likely due to the lower number  
320 of insect species compared to plant species modelled. Surprisingly, the most established  
321 thresholding methods for single species SDMs based on sensitivity and specificity (i.e.  
322 *Max.TSS*, *Opt.ROC* and *SenSpec*; Fig. 4 and Fig. S1 and S2) never ranked highest, as one or

323 more of the other thresholding method always ranked above them, both for community  
324 composition and for species richness.



## 325 **DISCUSSION**

### 326 **Do the most established thresholds for single species work as well for community** 327 **predictions?**

328 In this paper, we asked if the most established methods for single species thresholding are  
329 also the optimal choice for making predictions at the community level and if there is a direct  
330 link between the individual species predictions and the corresponding community metrics.  
331 Our results confirm the existence of such a link for single-index based metrics such as  
332 sensitivity, specificity and accuracy. However, these results should be interpreted with caution  
333 as maximising sensitivity or specificity can simply be achieved by predicting the species as  
334 present or absent (respectively) everywhere. In our study system, most of the modelled  
335 species have a low prevalence (i.e. are absent at most sites), thus accuracy (PCC) can often be  
336 improved by predicting the species as “absent” nearly everywhere.

337 The two most commonly used community evaluation metrics, Sørensen similarity index and  
338 deviation in species richness, were only weakly correlated with most evaluation metrics used  
339 for individual species. The most established thresholding methods for individual species  
340 predictions (i.e., *Max.TSS*, *Opt.ROC*, *SenSpec*) did show lower performance when applied to  
341 community-level predictions. This is likely due to the fact that both TSS and ROC try to find  
342 the best trade-off between sensitivity and specificity (Guisan, Thuiller & Zimmermann 2017).  
343 As most of the species have a prevalence far below 50% (i.e., are absent in many more sites  
344 than present), adding a few more presences might have a big effect on the sensitivity (by  
345 increasing the chance of finding the few real presences) but only marginally affects the  
346 specificity. By definition, increasing sensitivity also increases TSS, but with the drawback of  
347 a slight over-prediction. While this might not matter much on a single species basis, for  
348 community-level predictions the over-prediction will accumulate when summing binarised  
349 maps across all species, leading to the often observed over-estimation of species richness in S-

350 SDMs (e.g., Pineda & Lobo 2009; Dubuis *et al.* 2011; Mateo *et al.* 2012; Pottier *et al.* 2013;  
351 Pouteau *et al.* 2015; Zurell *et al.* 2016). It is important to remark, that in the rare case of an  
352 ecosystem mostly comprising of widespread species (i.e., prevalence >50 %) this will turn  
353 into the opposite as TSS and ROC will optimise absences leading to an underestimation of  
354 species richness. The strength of the over/under prediction bias is therefore linked to the  
355 prevalence distribution of the modelled species assemblages. However, in the vast majority of  
356 natural systems, both the site SR and the regional species pool are driven by a large number of  
357 rare (low prevalence species) compared to a few widespread species (Preston 1948; Magurran  
358 & Henderson 2003).

359 The community-based thresholding methods based on the selection of the most probable  
360 species (through a probability ranking) up to the predicted site richness (*MEM+PRR*, *pS-*  
361 *SDM+PRR*) can overcome this problem, because they are able to constrain species predictions  
362 based on a different value of species richness in each site (i.e. making them site-specific  
363 thresholding methods). Therefore, these methods prevent over-prediction while still allowing  
364 the analyses of species composition. Our results thus support the conclusion that, when the  
365 final goal is to optimize community composition, community-thresholding methods are the  
366 best option. Yet, as discussed in the next section, two single-species thresholding methods –  
367 *maximized Kappa* and *observed prevalence* – also showed good results for predicting  
368 communities (close to the community-based approaches). However, as community-based  
369 thresholds combine the optimisation of species richness prediction and a probability ranking  
370 rule (PRR), they would always select the species with the highest predicted probabilities in  
371 each site (D'Amen, Pradervand & Guisan 2015). This could seem logic and straightforward,  
372 but there might be a bias when the species in the community have varying prevalence  
373 (D'Amen *et al.* 2017a). In fact, the maximum predicted probability is depending on the  
374 prevalence of the species, thus the common species will tend to always have greater

375 maximum predicted probabilities than rare species and, as a result, will be considered present  
376 an over-proportionate number of time in the final community compositions. This bias will  
377 produce high similarity scores (Sørensen index) in the prediction evaluation, as the most  
378 common species are correctly predicted in most sites. However, the drawback is that the rarest  
379 species will be often omitted in the community predictions, which can be for instance  
380 problematic if the final goal of the modelling exercise has conservation implications.

### 381 **Is there a “best” threshold for community S-SDMs?**

382 We also tested if different methods for binarising community S-SDMs could be superior  
383 depending on the taxonomic group, prevalence distribution or species richness. While we  
384 observed significant differences between the different groups (i.e. taxa), there is no simple  
385 statistical way to assess if these differences are attributable to the biology of the taxa  
386 themselves or simply to the differences in site species richness and prevalence distributions.  
387 Nevertheless, when we standardized the deviation in species richness by the total number of  
388 modelled species (regional species pool), no significant difference was any more visible  
389 among the different taxonomic groups. The differences in species richness deviation seem  
390 therefore a direct cause of the regional species pool. The same also seems correct for the  
391 Sørensen similarity index, as datasets with higher species richness and species pool have  
392 lower similarity scores. This likely results from the fact that the more species need to be  
393 predicted correctly, the more difficult it becomes to predict the whole communities.

394 A similar ranking of thresholding methods was overall observed across taxonomic group  
395 within a given modelling pathway, while among the pathways there were clear shifts in the  
396 ranking of thresholding methods: with no independent community evaluation data (SSCV),  
397 the *Obs.Preval* and *Max.Kappa* threshold showed superior results, while the pathways using  
398 independent community evaluation data (ID and CCV) indicated the community-based  
399 thresholding to be superior (*pS-SDM+PRR*). This observation is in line with published

400 literature, where studies not using independent community data usually report a good  
401 performance of single species optimisations methods (e.g. D'Amen, Pradervand & Guisan  
402 2015; Distler *et al.* 2015; Thuiller *et al.* 2015), while studies using independent data usually  
403 have better results using community constraints (e.g. D'Amen *et al.* 2015). Yet, it is  
404 remarkable to notice that, although previously much criticized in the literature (e.g.,  
405 McPherson, Jetz & Rogers 2004; Allouche, Tsoar & Kadmon 2006), maximized Kappa  
406 (together here with the observed prevalence) did indeed perform well as a thresholding  
407 method for predicting both single species and communities, being nearly always superior to  
408 the sensitivity-specificity thresholding methods supporting earlier findings of Manel,  
409 Williams and Ormerod (2001).

410 It is important to notice that the shift in ranking between modelling pathways was likely due  
411 to a lower degree of overfitting and therefore a lower decrease in performance when  
412 predicting to independent data.

### 413 **Summing up: How to evaluate community predictions correctly?**

414 Our results show that the “single species cross-validation” approach (SSCV), the most  
415 commonly used in the literature to evaluate community predictions (e.g., Dubuis *et al.* 2011;  
416 Calabrese *et al.* 2014; Distler *et al.* 2015), yields overoptimistic and thus not fully realistic  
417 measures of predictive power. While this approach is usually able to provide satisfying  
418 evaluation for single species, as revealed by the cross-validation of individual species runs, it  
419 shows a clear degradation of predictions when measured at the level of communities. This  
420 occurs likely because “all” sites are used at least once at some stage across all modelling runs  
421 of the split-sampling procedure, and thus no observation (or very few in the best cases)  
422 remains fully independent (i.e. unused) for the final evaluation at the community level.  
423 Additionally, the sets of training sites used at each run differ among the species, making the  
424 results not entirely comparable across species.

425 The second approach found in the literature builds on the first one (SSCV; thus including an  
426 internal cross-validation evaluation), but uses spatially or temporally independent data (ID)  
427 for the assessment (thus an external evaluation), thus (unlike SSCV) using the same set of  
428 evaluation sites for all species (e.g., Benito, Cayuela & Albuquerque 2013; Pottier *et al.* 2013;  
429 Cord *et al.* 2014). When such independent data are available, this method provides the best  
430 possible evaluation, provided that the evaluation data are representative of the area where the  
431 models apply. This approach – with both internal and external evaluation - is also the one  
432 considered as optimal in James *et al.* (2013), and recently promoted in the field of SDMs by  
433 Guisan, Thuiller and Zimmermann (2017).

434 The third approach (CVV), newly presented here, repeats the ID approach a large number of  
435 times within a cross-validation procedure at the community-level (no example of this  
436 approach known in the literature). By doing this, the risk of bias in the evaluation data,  
437 inherent to the selection of a single evaluation data set, is minimized compared to the simple  
438 ID approach. Additionally, the repeated cross-validation allows the assessments of uncertainty  
439 and confidence intervals around the community predictions' performance metrics. However,  
440 as this approach selects the same sites for all species, its application is only possible under  
441 specific circumstances. First, all the species data need to be collected in the same sites (i.e.  
442 true 'community data'). Second, as this approach leads to an unequal number of  
443 presences/absences between different cross-validation runs for the same species, it can lead to  
444 models failing for very rare (low sample size) species in some of the cross-validation runs if  
445 not enough presence sites are selected in the training set.

446 According to our results and despite the potential limitations we advise the use of the  
447 proposed community cross-validation approach (CCV) to evaluate community models in  
448 future studies. In fact, we clearly showed that the common practice of evaluating the  
449 community predictions on the same dataset used for calibration process (SSCV) leads to

450 overoptimistic estimations of model performance. In the commonest case of unavailability of  
451 truly spatial (i.e., different region) or temporal (i.e., different sampling period) independent  
452 data, often independent datasets are “created” by randomly splitting the initial dataset in two  
453 parts. However, we advocate against this practise and instead promote the community cross-  
454 validation approach, which minimizes the artefacts of randomly splitting the initial data and  
455 allows the estimation of uncertainty associated with the community evaluation metrics.

456

## 457 **Acknowledgements**

458 This study was supported by the Swiss national Science Foundation (SESAM'ALP project,  
459 grant nr 31003A-1528661) to AG and by the European Commission, Marie Skłodowska-  
460 Curie Research Fellowship Programme (SESAM-ZOO project) to MDA and AG. R.G.M. was  
461 funded by a Marie Curie Intra-European Fellowship within the 7th European Community  
462 Framework Programme (ACONITE, PIEF-GA-2013-622620). The computations were  
463 performed at the Vital-IT (<http://www.vital-it.ch>) Center for high-performance computing of  
464 the SIB Swiss Institute of Bioinformatics.

## 465 **Authors' contributions**

466 DS and AG conceived the ideas; RF and MD analysed the plant and insect data; DS and RGM  
467 developed the modelling framework; DS led the writing and all authors contributed critically  
468 to the drafts and gave final approval for publication.

## 469 **Data Accessibility**

470 A generalised version of the community cross-validation algorithm is available in the ecospat  
471 R package (Cola *et al.* 2016) on GitHub (ecospat.CCV;  
472 <https://doi.org/10.5281/zenodo.1287805>). All species and environmental data are available on  
473 Dryad: <https://doi.org/10.5061/dryad.28d4k> (Grassland species and environmental predictors  
474 for plants; Guisan, Dubuis & Vittoz 2011) and <https://doi.org/10.5061/dryad.nf925ps> (forest,  
475 insect species and environmental predictors for insects; Guisan *et al.* 2018).

## 476 **REFERENCES**

477 Alexander, J.M., Diez, J.M., Hart, S.P. & Levine, J.M. (2016) When climate reshuffles competitors: a call for  
478 experimental macroecology. *Trends in Ecology & Evolution*, **31**, 831-841. doi:  
479 10.1016/j.tree.2016.08.003

480 Allouche, O., Tsoar, A. & Kadmon, R. (2006) Assessing the accuracy of species distribution models: prevalence,  
481 kappa and the true skill statistic (TSS). *Journal of Applied Ecology*, **43**, 1223-1232. doi:  
482 10.1111/j.1365-2664.2006.01214.x

483 Begon, M., Harper, J.L. & Townsend, C.R. (1996) *Ecology: individuals populations and communities*, Third  
484 edition edn. Blackwell Science Inc, Oxford.

485 Benito, B.M., Cayuela, L. & Albuquerque, F.S. (2013) The impact of modelling choices in the predictive  
486 performance of richness maps derived from species-distribution models: guidelines to build better  
487 diversitymodels. *Methods in Ecology and Evolution*, **4**, 327-335. doi: 10.1111/2041-210x.12022

488 Blois, J.L., Zarnetske, P.L., Fitzpatrick, M.C. & Finnegan, S. (2013) Climate Change and the Past, Present, and  
489 Future of Biotic Interactions. *Science*, **341**, 499-504. doi: 10.1126/science.1237184

490 Breiman, L. (2001) Random forests. *Machine Learning*, **45**, 5-32. doi: 10.1023/A:1010933404324

491 Breiner, F.T., Guisan, A., Bergamini, A. & Nobis, M.P. (2015) Overcoming limitations of modelling rare species  
492 by using ensembles of small models. *Methods in Ecology and Evolution*, **6**, 1210-1218. doi:  
493 10.1111/2041-210x.12403

494 Calabrese, J.M., Certain, G., Kraan, C. & Dormann, C.F. (2014) Stacking species distribution models and  
495 adjusting bias by linking them to macroecological models. *Global Ecology and Biogeography*, **23**, 99-  
496 112. doi: 10.1111/Geb.12102

497 Cohen, J. (1968) Weighted kappa: nominal scale agreement with provision for scaled disagreement or partial  
498 credit. *American Psychological Association*, **70**, 213-220.

499 Cola, V.D., Broennimann, O., Petitpierre, B., Breiner, F.T., D'Amen, M., Randin, C., . . . Dubuis, A. (2016)  
500 ecospat: an R package to support spatial analyses and modeling of species niches and distributions.  
501 *Ecography*.

502 Cord, A.F., Klein, D., Gernandt, D.S., la Rosa, J.A.P. & Dech, S. (2014) Remote sensing data can improve  
503 predictions of species richness by stacked species distribution models: a case study for Mexican pines.  
504 *Journal of Biogeography*, **41**, 736-748. doi: 10.1111/jbi.12225

505 Currie, D.J., Mittelbach, G.G., Cornell, H.V., Field, R., Guégan, J.F., Hawkins, B.A., . . . O'Brien, E. (2004)  
506 Predictions and tests of climate-based hypotheses of broad-scale variation in taxonomic richness.  
507 *Ecology Letters*, **7**, 1121-1134.



508 D'Amen, M., Bombi, P., Pearman, P.B., Schmatz, D.R., Zimmermann, N.E. & Bologna, M.A. (2011) Will  
509 climate change reduce the efficacy of protected areas for amphibian conservation in Italy? *Biological*  
510 *Conservation*, **144**, 989-997. doi: 10.1016/j.biocon.2010.11.004

511 D'Amen, M., Dubuis, A., Fernandes, R.F., Pottier, J., Pellissier, L. & Guisan, A. (2015) Using species richness  
512 and functional traits predictions to constrain assemblage predictions from stacked species distribution  
513 models. *Journal of Biogeography*, **42**, 1255-1266. doi: 10.1111/jbi.12485

514 D'Amen, M., Mateo, R.G., Pottier, J., Thuiller, W., Maiorano, L., Pellissier, L., . . . Guisan, A. (2017a)  
515 Improving spatial predictions of taxonomic, functional and phylogenetic diversity. *Journal of Ecology*.

516 D'Amen, M., Pradervand, J.N. & Guisan, A. (2015) Predicting richness and composition in mountain insect  
517 communities at high resolution: a new test of the SESAM framework. *Global Ecology and*  
518 *Biogeography*, **24**, 1443-1453. doi: 10.1111/geb.12357

519 D'Amen, M., Rahbek, C., Zimmermann, N.E. & Guisan, A. (2017b) Spatial predictions at the community level:  
520 from current approaches to future frameworks. *Biological Reviews*, **92**, 169–187. 10.1111/brv.12222

521 Distler, T., G., S.J., Velasquez-Tibata, J. & Langham, G.M. (2015) Stacked species distribution models and  
522 macroecological models provide congruent projections of avian species richness under climate change.  
523 *Journal of Biogeography*, **42**, 1-13. doi: 10.1111/jbi.12479

524 Dormann, C.F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carre, G., . . . Lautenbach, S. (2013) Collinearity: a  
525 review of methods to deal with it and a simulation study evaluating their performance. *Ecography*, **36**,  
526 27-46. doi: 10.1111/j.1600-0587.2012.07348.x

527 Dubuis, A., Pottier, J., Rion, V., Pellissier, L., Theurillat, J.P. & Guisan, A. (2011) Predicting spatial patterns of  
528 plant species richness: a comparison of direct macroecological and species stacking modelling  
529 approaches. *Diversity and Distributions*, **17**, 1122-1131. doi: 10.1111/j.1472-4642.2011.00792.x

530 Dullinger, S., Gattlinger, A., Thuiller, W., Moser, D., Zimmermann, N.E., Guisan, A., . . . Hulber, K. (2012)  
531 Extinction debt of high-mountain plants under twenty-first-century climate change. *Nature Climate*  
532 *Change*, **2**, 619-622. doi: 10.1038/Nclimate1514

533 Elith, J., Graham, C.H., Anderson, R.P., Dudik, M., Ferrier, S., Guisan, A., . . . Zimmermann, N.E. (2006) Novel  
534 methods improve prediction of species' distributions from occurrence data. *Ecography*, **29**, 129-151.  
535 doi: 10.1111/j.2006.0906-7590.04596.x

536 Elith, J. & Leathwick, J.R. (2009) Species Distribution Models: Ecological Explanation and Prediction Across  
537 Space and Time. *Annual Review of Ecology Evolution and Systematics*, **40**, 677-697. doi:  
538 10.1146/annurev.ecolsys.110308.120159

539 Elith, J., Leathwick, J.R. & Hastie, T. (2008) A working guide to boosted regression trees. *Journal of Animal*  
540 *Ecology*, **77**, 802-813. doi: 10.1111/j.1365-2656.2008.01390.x

541 Ferrier, S., Drielsma, M., Manion, G. & Watson, G. (2002) Extended statistical approaches to modelling spatial  
542 pattern in biodiversity in northeast New South Wales. II. Community-level modelling. *Biodiversity and*  
543 *Conservation*, **11**, 2309-2338. doi: 10.1023/A:1021374009951

544 Ferrier, S. & Guisan, A. (2006) Spatial modelling of biodiversity at the community level. *Journal of Applied*  
545 *Ecology*, **43**, 393-404. doi: 10.1111/j.1365-2664.2006.01149.x

546 Fielding, A.H. & Bell, J.F. (1997) A review of methods for the assessment of prediction errors in conservation  
547 presence-absence models. *Environmental Conservation*, **24**, 38-49. doi: 10.1017/S0376892997000088

548 Fleishman, E., Noss, R.F. & Noon, B.R. (2006) Utility and limitations of species richness metrics for  
549 conservation planning. *Ecological Indicators*, **6**, 543-553. doi: 10.1016/j.ecolind.2005.07.005

550 Freeman, E.A. & Moisen, G.G. (2008) A comparison of the performance of threshold criteria for binary  
551 classification in terms of predicted prevalence and kappa. *Ecological Modelling*, **217**, 48-58. doi:  
552 10.1016/j.ecolmodel.2008.05.015

553 Gotelli, N.J., Anderson, M.J., Arita, H.T., Chao, A., Colwell, R.K., Currie, D.J., . . . Willig, M.R. (2009) Patterns  
554 and causes of species richness: a general simulation model for macroecology. *Ecology Letters*, **12**, 873–  
555 886. doi: 10.1111/j.1461-0248.2009.01353.x

556 Granger, V., Bez, N., Fromentin, J.M., Meynard, C., Jadaud, A. & Merigot, B. (2015) Mapping diversity indices:  
557 not a trivial issue. *Methods in Ecology and Evolution*, **6**, 688-696.

558 Guisan, A., Dubuis, A., Pellisser, L., Pradervand, J.N., Meier, S., Scherrer, D., . . . Vittoz, P. (2018) Data from:  
559 How to best threshold and validate stacked species assemblages? Community optimisation might hold  
560 the answer. *Dryad Digital Repository*. <https://doi.org/10.5061/dryad.nf925ps>

561 Guisan, A., Dubuis, A. & Vittoz, P. (2011) Data from: Predicting spatial patterns of plant species richness:  
562 a comparison of direct macroecological and species stacking modelling approaches. *Dryad Digital*  
563 *Repository*. <https://doi.org/10.5061/dryad.28d4k>

564 Guisan, A., Lehmann, A., Ferrier, S., Aspinall, R., Overton, R., Austin, M. & Hastie, T. (2006) Making better  
565 biogeographic predictions of species distribution. *Journal of Applied Ecology*, **43**, 386-392.

566 Guisan, A. & Rahbek, C. (2011) SESAM - a new framework integrating macroecological and species  
567 distribution models for predicting spatio-temporal patterns of species assemblages. *Journal of*  
568 *Biogeography*, **38**, 1433-1444. doi: 10.1111/j.1365-2699.2011.02550.x

569 Guisan, A. & Thuiller, W. (2005) Predicting species distribution: offering more than simple habitat models.  
570 *Ecology Letters*, **8**, 993-1009. doi: 10.1111/j.1461-0248.2005.00792.x

571 Guisan, A., Thuiller, W. & Zimmermann, N.E. (2017) *Habitat suitability and distribution models*. Cambridge  
572 University Press.

573 Guisan, A., Tingley, R., Baumgartner, J.B., Naujokaitis-Lewis, I., Sutcliffe, P.R., Tulloch, A.I.T., . . . Buckley,  
574 Y.M. (2013) Predicting species distributions for conservation decisions. *Ecology Letters*, **16**, 1424-  
575 1435. doi: 10.1111/Ele.12189

576 Hartmann, P., Fouvry, P. & Horisberger, D. (2009) L'Observatoire de l'écosystème forestier du canton de Vaud:  
577 espace de recherche appliquée| The Forest Ecosystem Observatory in Canton Vaud: a field of applied  
578 research. *Schweizerische Zeitschrift für Forstwesen*, **160**, s2-s6.

579 Hastie, T.J. & Tibshirani, R. (1990) *Generalized Additive Models*. Chapman & Hall, London.

580 Hawkins, B.A. & Porter, E.E. (2003) Water–energy balance and the geographic pattern of species richness of  
581 western Palearctic butterflies. *Ecological Entomology*, **28**, 678-686. doi: 10.1111/j.1365-  
582 2311.2003.00551.x

583 Hespanhol, H., Cezón, K., Felicísimo, Á.M., Muñoz, J. & Mateo, R.G. (2015) How to describe species richness  
584 patterns for bryophyte conservation? *Ecology and evolution*, **5**, 5443-5455.

585 James, G., Witten, D., Hastie, T. & Tibshirani, R. (2013) *An introduction to statistical learning*. Springer.

586 Jimenez-Valverde, A. & Lobo, J.M. (2007) Threshold criteria for conversion of probability of species presence  
587 to either-or presence-absence. *Acta Oecologica-International Journal of Ecology*, **31**, 361-369. doi:  
588 10.1016/j.actao.2007.02.001

589 Kearney, M. & Porter, W. (2009) Mechanistic niche modelling: combining physiological and spatial data to  
590 predict species' ranges. *Ecology Letters*, **12**, 334-350. doi: 10.1111/j.1461-0248.2008.01277.x

591 Liu, C., White, M. & Newell, G. (2013) Selecting thresholds for the prediction of species occurrence with  
592 presence-only data. *Journal of Biogeography*, **40**, 778-789.

593 Liu, C.R., Berry, P.M., Dawson, T.P. & Pearson, R.G. (2005) Selecting thresholds of occurrence in the  
594 prediction of species distributions. *Ecography*, **28**, 385-393. doi: 10.1111/j.0906-7590.2005.03957.x

595 Lomba, A., Pellissier, L., Randin, C., Vicente, J., Moreira, F., Honrado, J. & Guisan, A. (2010) Overcoming the  
596 rare species modelling paradox: A novel hierarchical framework applied to an Iberian endemic plant.  
597 *Biological Conservation*, **143**, 2647-2657. doi: 10.1016/j.biocon.2010.07.007

598 Maggini, R., Lehmann, A., Zimmermann, N.E. & Guisan, A. (2006) Improving generalized regression analysis  
599 for the spatial prediction of forest communities. *Journal of Biogeography*, **33**, 1729-1749. doi:  
600 10.1111/j.1365-2699.2006.01465.x

601 Magurran, A.E. & Henderson, P.A. (2003) Explaining the excess of rare species in natural species abundance  
602 distributions. *Nature*, **422**, 714-716.

603 Manel, S., Williams, H.C. & Ormerod, S.J. (2001) Evaluating presence-absence models in ecology: the need to  
604 account for prevalence. *Journal of Applied Ecology*, **38**, 921-931. doi: 10.1046/j.1365-  
605 2664.2001.00647.x

606 Marmion, M., Parviainen, M., Luoto, M., Heikkinen, R.K. & Thuiller, W. (2009) Evaluation of consensus  
607 methods in predictive species distribution modelling. *Diversity and Distributions*, **15**, 59-69. doi:  
608 10.1111/j.1472-4642.2008.00491.x

609 Mateo, R.G., de la Estrella, M., Felicísimo, A.M., Muñoz, J. & Guisan, A. (2013) A new spin on a  
610 compositionalist predictive modelling framework for conservation planning: A tropical case study in  
611 Ecuador. *Biological Conservation*, **160**, 150-161. doi: 10.1016/j.biocon.2013.01.014

612 Mateo, R.G., Felicísimo, A.M., Pottier, J., Guisan, A. & Muñoz, J. (2012) Do Stacked Species Distribution  
613 Models Reflect Altitudinal Diversity Patterns? *PLoS ONE*, **7**, e32586. doi:  
614 10.1371/journal.pone.0032586

615 Mateo, R.G., Mokany, K. & Guisan, A. (2017) Biodiversity Models: What If Unsaturation Is the Rule? *Trends in*  
616 *Ecology & Evolution*, **32**, 556-566. doi: <http://dx.doi.org/10.1016/j.tree.2017.05.003>

617 McCullagh, P. & Nelder, J.A. (1989) *Generalized Linear Models. 2nd edition*. Chapman and Hall, London.

618 McPherson, J.M., Jetz, W. & Rogers, D.J. (2004) The effects of species' range sizes on the accuracy of  
619 distribution models: ecological phenomenon or statistical artefact? *Journal of Applied Ecology*, **41**, 811-  
620 823. doi: 10.1111/j.0021-8901.2004.00943.x

621 Meier, E.S., Kienast, F., Pearman, P.B., Svenning, J.C., Thuiller, W., Araujo, M.B., . . . Zimmermann, N.E.  
622 (2010) Biotic and abiotic variables show little redundancy in explaining tree species distributions.  
623 *Ecography*, **33**, 1038-1048. doi: 10.1111/j.1600-0587.2010.06229.x

624 Merow, C., Smith, M.J., Edwards, T.C., Guisan, A., McMahon, S.M., Normand, S., . . . Elith, J. (2014) What do  
625 we gain from simplicity versus complexity in species distribution models? *Ecography*, **37**, 1267-1281.  
626 doi: 10.1111/ecog.00845

627 Mod, H.K., le Roux, P.C., Guisan, A. & Luoto, M. (2015) Biotic interactions boost spatial models of species  
628 richness. *Ecography*, **38**, 913-921. doi: 10.1111/ecog.01129

629 Mokany, K. & Ferrier, S. (2011) Predicting impacts of climate change on biodiversity: a role for semi-  
630 mechanistic community-level modelling. *Diversity and Distributions*, **17**, 374-380. doi: 10.1111/j.1472-  
631 4642.2010.00735.x

632 Nenzen, H.K. & Araujo, M.B. (2011) Choice of threshold alters projections of species range shifts under climate  
633 change. *Ecological Modelling*, **222**, 3346-3354. doi: 10.1016/j.ecolmodel.2011.07.011

634 Nogues-Bravo, D. & Rahbek, C. (2011) Communities Under Climate Change. *Science*, **334**, 1070-1071. doi:  
635 10.1126/science.1214833

636 Parmesan, C. & Yohe, G. (2003) A globally coherent fingerprint of climate change impacts across natural  
637 systems. *Nature*, **421**, 37-42. doi: 10.1038/nature01286

638 Pellissier, L., Pradervand, J.-N., Pottier, J., Dubuis, A., Maiorano, L. & Guisan, A. (2012) Climate-based  
639 empirical models show biased predictions of butterfly communities along environmental gradients.  
640 *Ecography*, **35**, 684-692. doi: 10.1111/j.1600-0587.2011.07047.x

641 Pereira, H.M., Ferrier, S., Walters, M., Geller, G.N., Jongman, R.H.G., Scholes, R.J., . . . Wegmann, M. (2013)  
642 Essential Biodiversity Variables. *Science*, **339**, 277-278. doi: 10.1126/science.1229931

643 Pineda, E. & Lobo, J.M. (2009) Assessing the accuracy of species distribution models to predict amphibian  
644 species richness patterns. *Journal of Animal Ecology*, **78**, 182-190. doi: 10.1111/j.1365-  
645 2656.2008.01471.x

646 Pottier, J., Dubuis, A., Pellissier, L., Maiorano, L., Rossier, L., Randin, C.F., . . . Guisan, A. (2013) The accuracy  
647 of plant assemblage prediction from species distribution models varies along environmental gradients.  
648 *Global Ecology and Biogeography*, **22**, 52-63. doi: 10.1111/j.1466-8238.2012.00790.x

649 Pouteau, R., Bayle, E., Blanchard, E., Birnbaum, P., Cassan, J.J., Hequet, V., . . . Vandrot, H. (2015) Accounting  
650 for the indirect area effect in stacked species distribution models to map species richness in a montane  
651 biodiversity hotspot. *Diversity and Distributions*, **21**, 1329-1338. doi: 10.1111/ddi.12374

652 Pradervand, J.N., Dubuis, A., Reymond, A., Sonnay, V., Gelin, A. & Guisan, A. (2013) Quels facteurs  
653 influencent la richesse en orthoptères des Préalpes vaudoises? *Bulletin de la Société Vaudoises des*  
654 *Sciences Naturelles*, **93**, 155-173.

655 Preston, F.W. (1948) The Commonness, and Rarity, of Species. *Ecology*, **29**, 254-283.

656 Scherrer, D., Massy, S., Meier, S., Vittoz, P. & Guisan, A. (2017) Assessing and predicting shifts in mountain  
657 forest composition across 25 years of climate change. *Diversity and Distributions*, **23**, 517-528.

658 Sørensen, T. (1948) A method of establishing groups of equal amplitude in plant sociology based on similarity of  
659 species and its application to analyses of the vegetation on Danish commons. *Biol. Skr.*, **5**, 1-34.

660 Thuiller, W., Georges, D. & Engler, R. (2013) biomod2: Ensemble platform for species distribution modeling. *R*  
661 *package version*, **2**, r560.

662 Thuiller, W., Lavorel, S., Araújo, M.B., Sykes, M.T. & Prentice, I.C. (2005) Climate change threats to plant  
663 diversity in Europe. *Proceedings of the National Academy of Sciences of the United States of America*,  
664 **102**, 8245-8250. doi: 10.1073/pnas.0409902102

665 Thuiller, W., Pollock, L.J., Gueguen, M. & Münkemüller, T. (2015) From species distributions to meta-  
666 communities. *Ecology Letters*, **18**, 1321-1328. doi: 10.1111/ele.12526

667 Turner, J.R., Gatehouse, C.M. & Corey, C.A. (1987) Does solar energy control organic diversity? Butterflies,  
668 moths and the British climate. *Oikos*, 195-205.

669 Van der Putten, W.H., Macel, M. & Visser, M.E. (2010) Predicting species distribution and abundance responses  
670 to climate change: why it is essential to include biotic interactions across trophic levels. *Philosophical*  
671 *Transactions of the Royal Society B-Biological Sciences*, **365**, 2025-2034. doi: 10.1098/rstb.2010.0037

672 Zimmermann, N.E., Edwards, T.C., Graham, C.H., Pearman, P.B. & Svenning, J.-C. (2010) New trends in  
673 species distribution modelling. *Ecography*, **33**, 985-989. doi: 10.1111/j.1600-0587.2010.06953.x

674 Zurell, D., Zimmermann, N.E., Sattler, T., Nobis, M.P. & Schröder, B. (2016) Effects of functional traits on the  
675 prediction accuracy of species richness models. *Diversity and Distributions*.

676

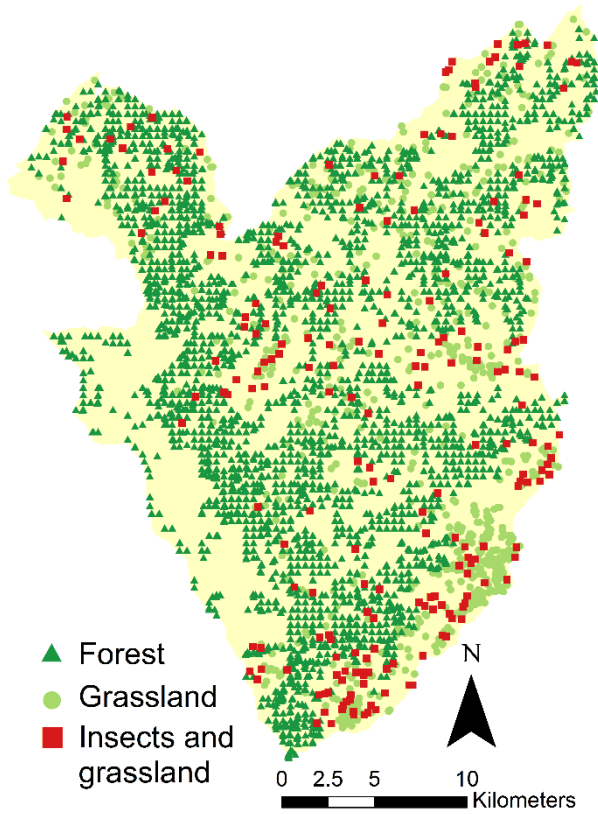
677 **Figure legends**

678 **Figure 1:** Map of the study area with the forested sites (dark green triangles, N=3076), the  
679 grassland sites (light green circles and red squares, N=903) and the insect sites (red squares,  
680 butterflies N=192, grasshoppers N=202).

681 **Figure 2:** The modelling framework illustrating the three different community modelling  
682 approaches: “single species cross-validation” (SSCV), “independent data” (ID) and  
683 “community cross-validation” (CCV).

684 **Figure 3:** Deviation in site specific species richness between observations and predictions for  
685 the four different datasets (top to bottom) and the three different modelling pathways (left to  
686 right). The boxplots are sorted by the median and the colours indicate the different  
687 thresholding techniques used to binarise predictions. The line in the box indicates the median,  
688 boxes range from the 25<sup>th</sup> to the 75<sup>th</sup> percentile and the whiskers indicate  $\pm 2$  standard  
689 deviations. Letters above the boxplots indicate significant differences (Wilcoxon rank sum  
690 test,  $p < 0.05$ ).

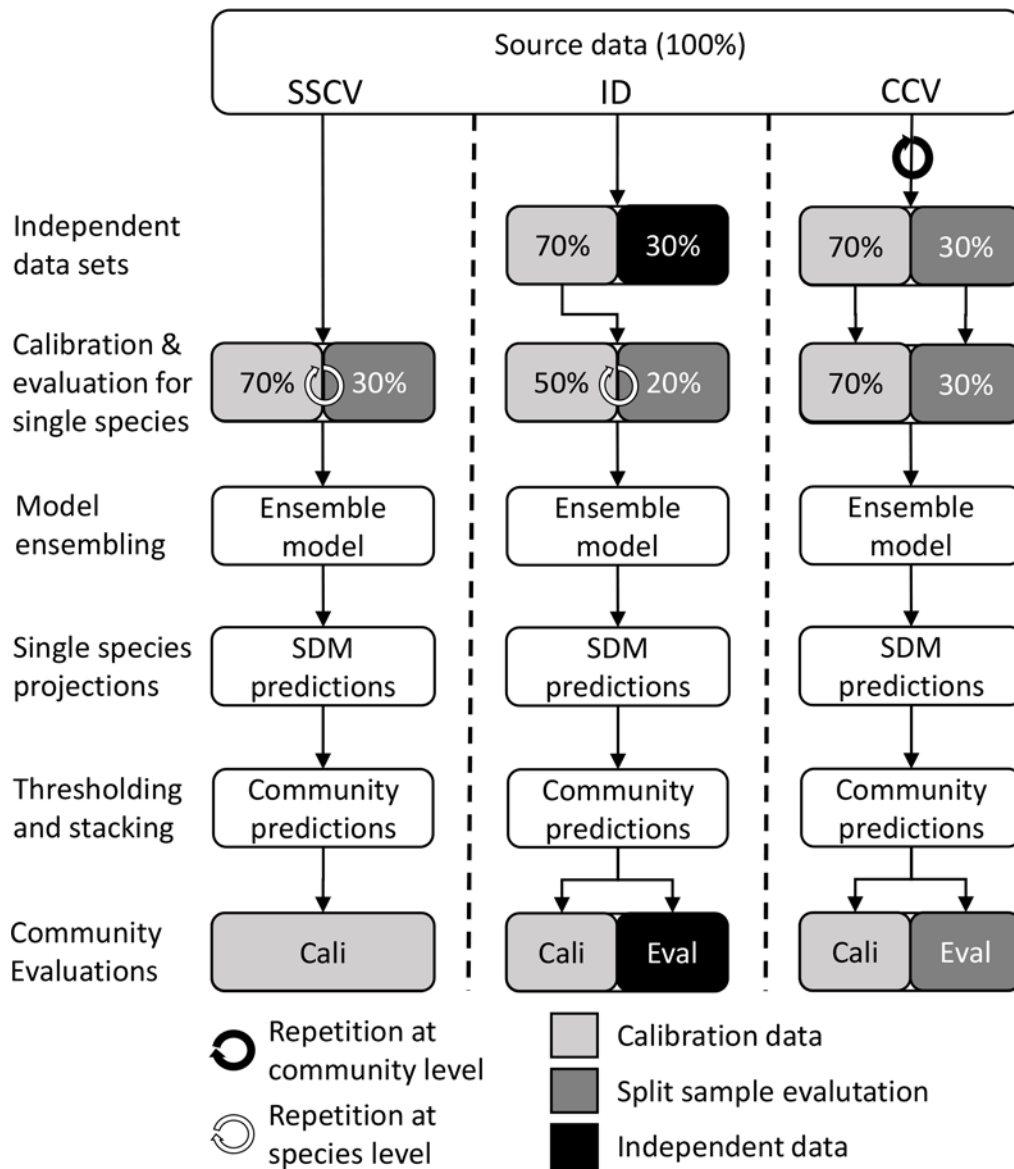
691 **Figure 4:** Sørensen similarity between observations and predictions for the four different  
692 datasets (top to bottom) and the three different modelling pathways (left to right). The  
693 boxplots are sorted by the median and the colours indicate the different thresholding  
694 techniques. The line in the box indicates the median, boxes range from the 25<sup>th</sup> to the 75<sup>th</sup>  
695 percentile and the whiskers indicate  $\pm 2$  standard deviations. Letters above the boxplots  
696 indicate significant differences (Wilcoxon rank sum test,  $p < 0.05$ ).



698

699 **Figure 1:** Map of the study area with the forested sites (dark green triangles, N=3076), the  
700 grassland sites (light green circles and red squares, N=903) and the insect sites (red squares,  
701 butterflies N=192, grasshoppers N=202).



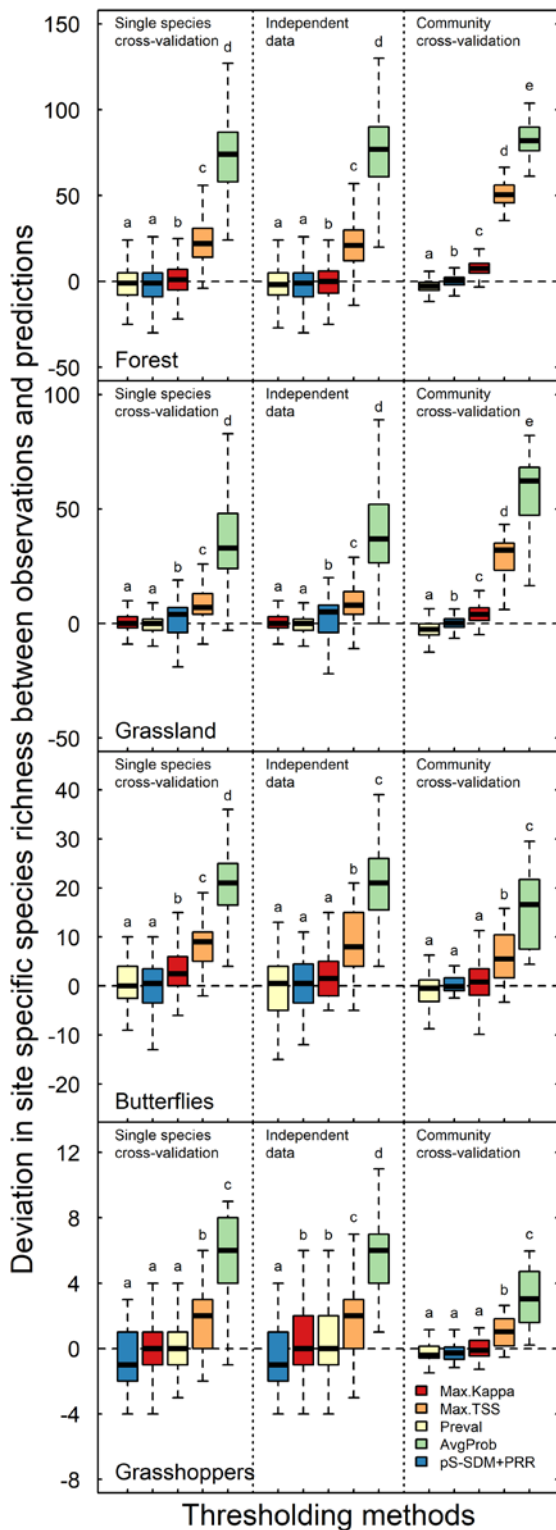


703

704 **Figure 2:** The modelling framework illustrating the three different community modelling

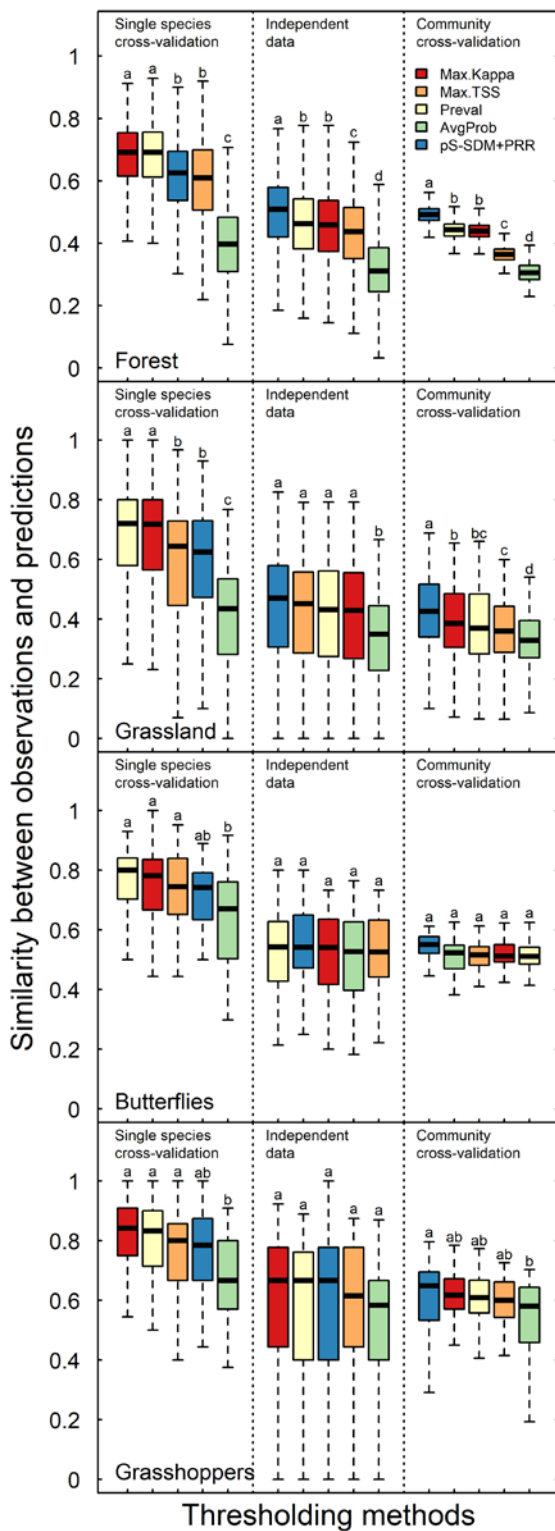
705 approaches: “single species cross-validation” (SSCV), “independent data” (ID) and

706 “community cross-validation” (CCV).



**Figure 3:** Deviation in site specific species richness between observations and predictions for the four different datasets (top to bottom) and the three different modelling pathways (left to right). The boxplots are sorted by the median and the colours indicate the different thresholding techniques used to binarise predictions. The line in the box indicates the median, boxes range from the 25<sup>th</sup> to the 75<sup>th</sup> percentile and the whiskers indicate  $\pm 2$  standard deviations. Letters above the boxplots indicate significant differences (Wilcoxon rank sum test,  $p < 0.05$ ).

**Figure 4**



**Figure 4:** Sørensen similarity between observations and predictions for the four different datasets (top to bottom) and the three different modelling pathways (left to right). The boxplots are sorted by the median and the colours indicate the different thresholding techniques. The line in the box indicates the median, boxes range from the 25<sup>th</sup> to the 75<sup>th</sup> percentile and the whiskers indicate  $\pm 2$  standard deviations. Letters above the boxplots indicate significant differences (Wilcoxon rank sum test,  $p < 0.05$ ).

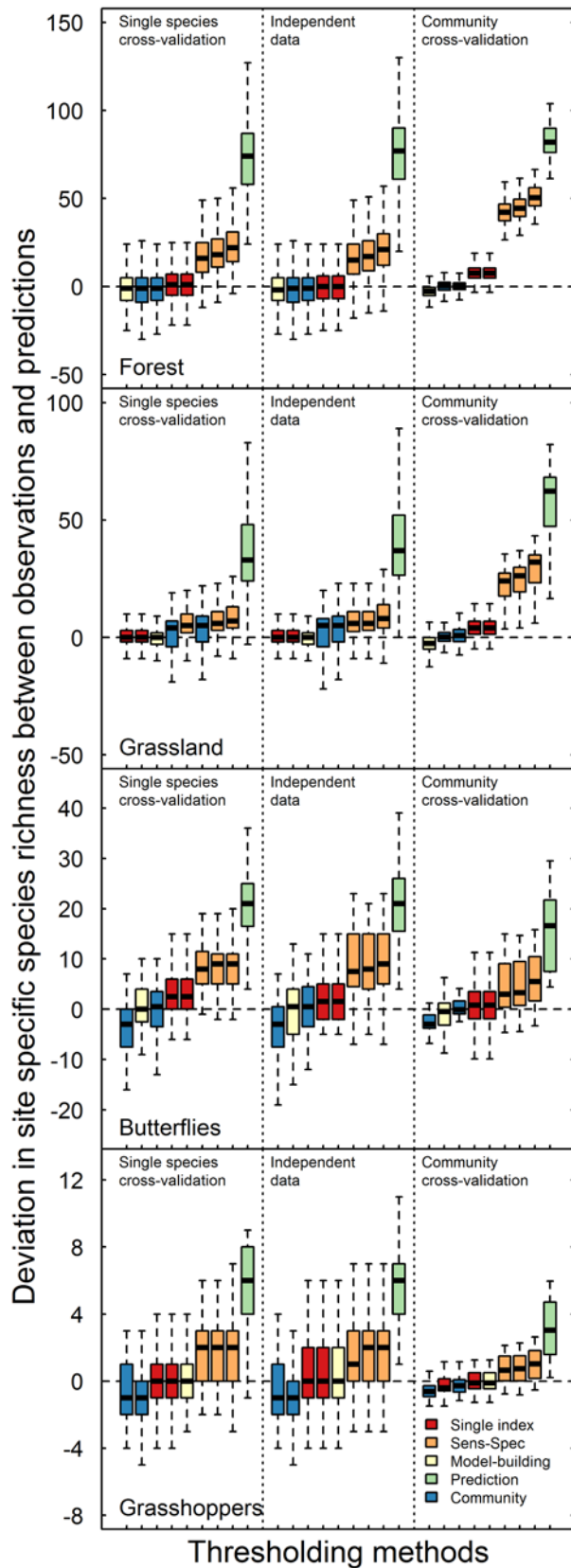
**Table 1:** Basic statistics of the data sets used for the case study and the evaluation metrics (AUC) for the individual species distribution models using the three different community evaluation approaches. SSCV = Single species cross-validation, ID = Independent data, CCV = Community cross-validation

Data set	Number of species modelled (recorded)	Prevalence (mean $\pm$ sd)	Species richness (mean $\pm$ sd)	AUC SSCV (mean $\pm$ sd)	AUC ID (mean $\pm$ sd)	AUC CCV (mean $\pm$ sd)
Forest	312 (703)	0.044 $\pm$ 0.090	29.5 $\pm$ 11.8	0.80 $\pm$ 0.09	0.80 $\pm$ 0.08	0.79 $\pm$ 0.09
Grassland	212 (905)	0.098 $\pm$ 0.089	23.5 $\pm$ 13.8	0.82 $\pm$ 0.07	0.83 $\pm$ 0.06	0.81 $\pm$ 0.06
Butterflies	77 (131)	0.235 $\pm$ 0.137	18.1 $\pm$ 9.2	0.76 $\pm$ 0.10	0.75 $\pm$ 0.12	0.76 $\pm$ 0.10
Grasshoppers	20 (41)	0.256 $\pm$ 0.193	5.1 $\pm$ 3.3	0.84 $\pm$ 0.07	0.86 $\pm$ 0.08	0.84 $\pm$ 0.06

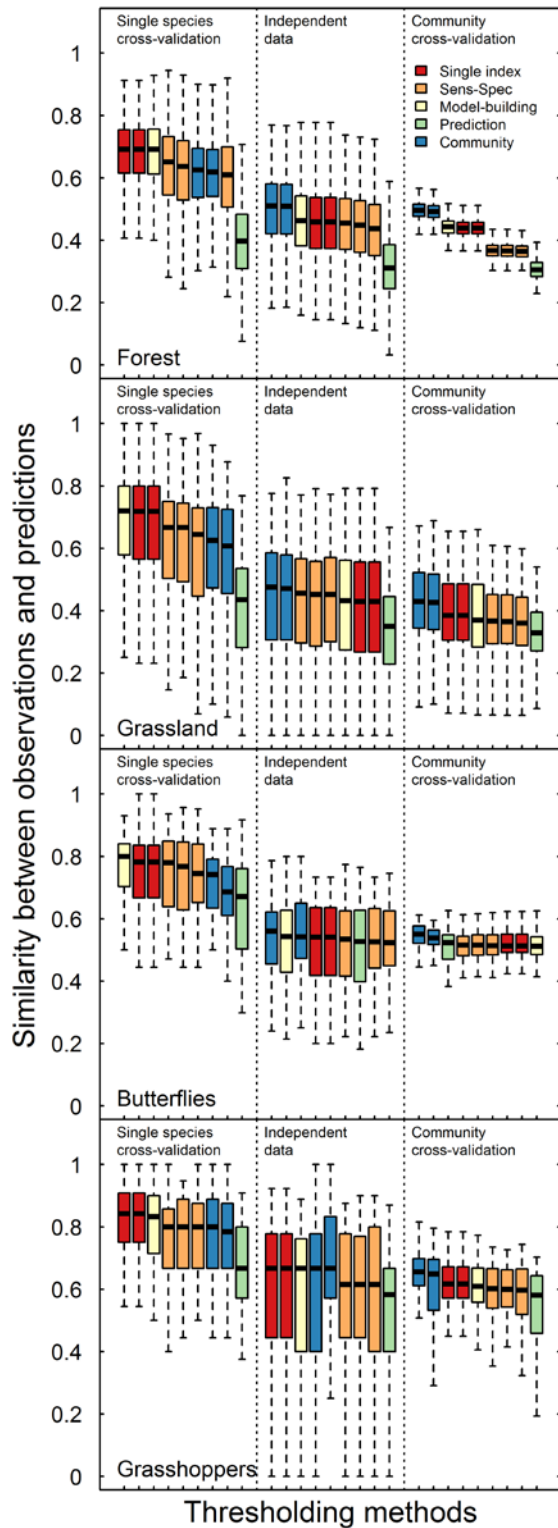
**Table 2:** Pearson Correlation of single species and community evaluation statistics. The asterisks indicate the significance level. Correlations of the single species evaluation metrics and their corresponding community evaluation metric are highlighted in bold.

Single species	Community metrics						
	Accuracy	Sensitivity	Specificity	KAPPA	TSS	Sørensen similarity	SR deviation
<b>Accuracy</b>	<b>1.00</b> ***	-0.37 *	0.95 ***	0.70 ***	0.37 *	0.37 *	-0.58 ***
<b>Sensitivity</b>	-0.36 **	<b>0.93</b> ***	-0.54 ***	0.01 n.s.	0.56 ***	0.18 n.s.	-0.44 ***
<b>Specificity</b>	0.97 ***	-0.53 ***	<b>0.99</b> ***	0.64 ***	0.20 n.s.	0.31 *	-0.63 ***
<b>KAPPA</b>	0.41 **	0.50 *	0.27 *	<b>0.79</b> ***	0.72 ***	0.82 ***	-0.3 *
<b>TSS</b>	0.06 n.s.	0.85 ***	-0.14 n.s.	0.35 n.s.	<b>0.79</b> ***	0.38 **	-0.20 n.s.

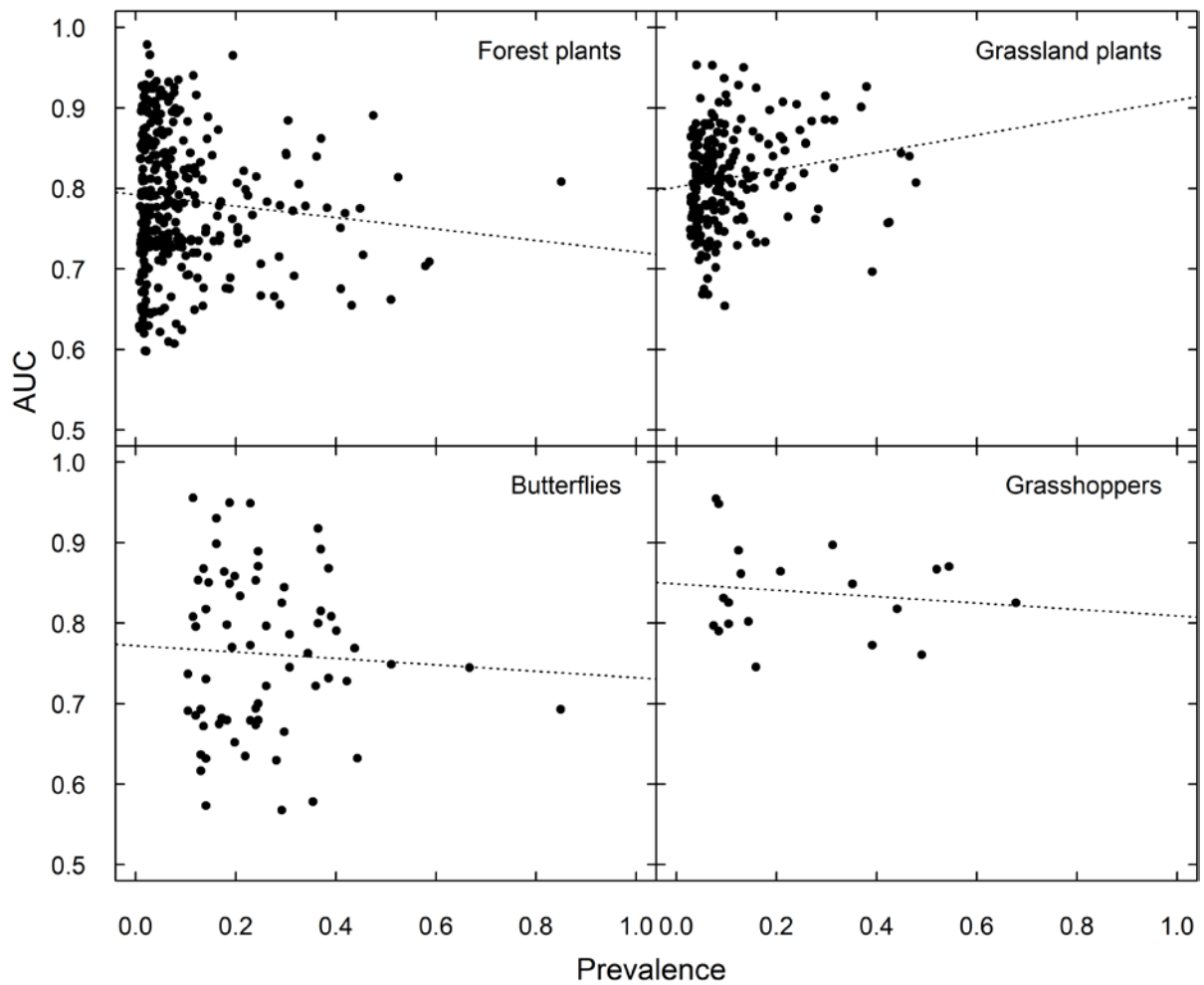
The asterisks indicate the significance level (n.s.= not significant, \* p<0.05, \*\* p<0.01, \*\*\* p<0.001)



**Figure S1:** Deviation in site specific species richness between observations and predictions for the four different datasets (top to bottom) and the three different modelling pathways (left to right). The boxplots are sorted by the median and the colours indicate the different thresholding techniques. The line in the box indicates the median, boxes range from the 25<sup>th</sup> to the 75<sup>th</sup> percentile and the whiskers indicate  $\pm 2$  standard deviations. For details on the method used within each threshold group see Table S1.



**Figure S2:** Sørensen similarity between observations and predictions for the four different datasets (top to bottom) and the three different modelling pathways (left to right). The boxplots are sorted by the median and the colours indicate the different thresholding techniques. The line in the box indicates the median, boxes range from the 25<sup>th</sup> to the 75<sup>th</sup> percentile and the whiskers indicate  $\pm 2$  standard deviations. For details on the method used within each threshold group see Table S1.



**Figure S3:** The relationship of the prevalence of a species (i.e., percentage of sites inhabited) to the performance of the SDMs (i.e., as measured by AUC) for the four studied data sets (taxa).



**Table S1:** Description of the ten thresholding methods based on Liu *et al.* (2005) and Nenzen and Araujo (2011).

Approach	Acronym	Definition	Reference
<b>Single index-based approaches</b>			
1. Kappa maximization approach	Max.Kappa	Kappa statistic is maximized	(Huntley <i>et al.</i> 1995; Guisan, Theurillat & Kienast 1998)
2. Maximum commission error	MCE05	Allowed a maximum commission error of 5%	(Mateo <i>et al.</i> 2012)
<b>Sensitivity and specificity-combined approaches</b>			
3. TSS maximization approach	Max.TSS	TSS statistic is maximized	(Allouche, Tsoar & Kadmon 2006)
4. Sensitivity-specificity equality approach	SensSpec	Difference of sens-spec is minimized	(Cantor <i>et al.</i> 1999)
5. ROC plot-based approach	Opt.ROC	ROC statistic is maximized	(Cantor <i>et al.</i> 1999)
<b>Model-building data-only-based approach</b>			
6. Prevalence approach	Preval	Prevalence of the calibration data	(Cramer 2003)
<b>Predicted probability-based approaches</b>			
7. Average probability approach	AvgProb	Taking the average predicted probability of the model-building data as threshold	(Cramer 2003)
<b>Community based approaches</b>			
8. pS-SDM+PRR	pS-SDM+PRR	Probability stacked SDM	(Dubuis <i>et al.</i> 2013)
9. MEM+PRR	MEM+PRR	Macroecological model for SR	(Guisan & Rahbek 2011)

**Table S2** : Community evaluation metrics used in this study.

Metric	Definition
<b>Species richness</b>	
Deviation in species richness	$Dev.SPR = n_{pred} - n_{obs}$
<b>Prediction success</b>	
Sensitivity	$Sens = \frac{TP}{TP + FA}$
Specificity	$Spec = \frac{TA}{TA + FP}$
Community accuracy	$Acc = \frac{TP + TA}{N}$
Community TSS	$TSS = Sens + Spec - 1$
Community Kappa	$K = \frac{Acc - p_e}{1 - p_e}$
<b>Community composition</b>	
Sørensen	$S = \frac{2 * TP}{2 * TP + FP + FA}$

$n_{pred}$  = Number of species predicted

$n_{obs}$  = Number of species observed

$N$  = Number of events

$TP$  = Correctly predicted present species

$TA$  = Correctly predicted absent species

$FP$  = Falsely predicted present species

$FA$  = Falsely predicted absent species

$p_e = \frac{(TP+FA)(TP+FP)+(TA+FP)(TA+FA)}{N^2}$

2 **Table S3** : Evaluation scores of individual SDMs by TSS (A), Kappa (B), PCC (C), Sensitivity (D) and Specificity (E) for the three community  
3 evaluation approaches and four datasets. SSCV = Single species cross-validation, ID = Independent data, CCV = Community cross-validation,  
4 FO = Forest plants, GL = Grassland plants, BF = Butterflies, GH = Grasshoppers.

5 (A) TSS

Thresholding Approach	SSCV				ID				CCV			
	FO	GL	BF	GH	FO	GL	BF	GH	FO	GL	BF	GH
Max.Kappa	0.2 ± 0.14	0.27 ± 0.2	0.31 ± 0.23	0.42 ± 0.17	0.21 ± 0.17	0.25 ± 0.21	0.3 ± 0.22	0.43 ± 0.26	0.23 ± 0.13	0.28 ± 0.14	0.31 ± 0.18	0.37 ± 0.15
MCE05	0.3 ± 0.17	0.32 ± 0.17	0.27 ± 0.18	0.43 ± 0.16	0.28 ± 0.17	0.34 ± 0.17	0.27 ± 0.23	0.45 ± 0.21	0.25 ± 0.17	0.34 ± 0.12	0.31 ± 0.17	0.42 ± 0.12
Max.TSS	0.35 ± 0.2	0.38 ± 0.14	0.34 ± 0.24	0.47 ± 0.12	0.33 ± 0.22	0.38 ± 0.21	0.34 ± 0.23	0.5 ± 0.26	0.35 ± 0.14	0.39 ± 0.11	0.35 ± 0.18	0.44 ± 0.12
SensSpec	0.32 ± 0.16	0.36 ± 0.14	0.33 ± 0.18	0.51 ± 0.18	0.31 ± 0.2	0.37 ± 0.21	0.34 ± 0.24	0.51 ± 0.26	0.35 ± 0.13	0.38 ± 0.11	0.35 ± 0.18	0.45 ± 0.11
Opt.ROC	0.34 ± 0.21	0.36 ± 0.19	0.33 ± 0.17	0.44 ± 0.26	0.32 ± 0.21	0.37 ± 0.21	0.34 ± 0.23	0.49 ± 0.25	0.35 ± 0.13	0.38 ± 0.11	0.35 ± 0.18	0.44 ± 0.12
Preval	0.18 ± 0.15	0.27 ± 0.15	0.3 ± 0.19	0.41 ± 0.23	0.2 ± 0.16	0.26 ± 0.2	0.31 ± 0.22	0.4 ± 0.23	0.18 ± 0.14	0.23 ± 0.15	0.3 ± 0.17	0.37 ± 0.17
AvgProb	0.43 ± 0.16	0.5 ± 0.12	0.41 ± 0.21	0.55 ± 0.14	0.47 ± 0.16	0.53 ± 0.13	0.38 ± 0.23	0.56 ± 0.15	0.44 ± 0.16	0.49 ± 0.11	0.4 ± 0.18	0.54 ± 0.14
pS-SDM+PRR	0.12 ± 0.19	0.17 ± 0.24	0.28 ± 0.24	0.28 ± 0.22	0.14 ± 0.18	0.2 ± 0.24	0.24 ± 0.24	0.29 ± 0.28	0.14 ± 0.17	0.19 ± 0.21	0.27 ± 0.23	0.27 ± 0.22
MEM+PRR	0.16 ± 0.17	0.2 ± 0.23	0.25 ± 0.24	0.3 ± 0.24	0.14 ± 0.18	0.2 ± 0.24	0.22 ± 0.24	0.32 ± 0.28	0.14 ± 0.17	0.2 ± 0.22	0.25 ± 0.22	0.3 ± 0.22

6

7 (B) KAPPA

Thresholding Approach	SSCV				ID				CCV			
	FO	GL	BF	GH	FO	GL	BF	GH	FO	GL	BF	GH
Max.Kappa	0.2 ± 0.14	0.24 ± 0.18	0.28 ± 0.22	0.35 ± 0.19	0.21 ± 0.15	0.24 ± 0.19	0.29 ± 0.22	0.42 ± 0.25	0.2 ± 0.14	0.24 ± 0.15	0.29 ± 0.18	0.36 ± 0.15
MCE05	0.21 ± 0.13	0.23 ± 0.15	0.24 ± 0.22	0.31 ± 0.16	0.21 ± 0.13	0.28 ± 0.15	0.27 ± 0.22	0.41 ± 0.18	0.11 ± 0.12	0.16 ± 0.14	0.24 ± 0.17	0.32 ± 0.16
Max.TSS	0.19 ± 0.12	0.27 ± 0.17	0.3 ± 0.19	0.4 ± 0.14	0.21 ± 0.14	0.27 ± 0.16	0.3 ± 0.21	0.41 ± 0.23	0.17 ± 0.13	0.22 ± 0.15	0.3 ± 0.17	0.36 ± 0.15
SensSpec	0.21 ± 0.14	0.21 ± 0.17	0.29 ± 0.19	0.41 ± 0.17	0.22 ± 0.14	0.28 ± 0.17	0.3 ± 0.21	0.42 ± 0.23	0.17 ± 0.13	0.23 ± 0.15	0.3 ± 0.17	0.37 ± 0.14
Opt.ROC	0.15 ± 0.15	0.22 ± 0.14	0.32 ± 0.18	0.43 ± 0.17	0.22 ± 0.14	0.27 ± 0.16	0.3 ± 0.21	0.42 ± 0.23	0.17 ± 0.13	0.23 ± 0.15	0.3 ± 0.18	0.37 ± 0.14
Preval	0.2 ± 0.13	0.22 ± 0.18	0.3 ± 0.2	0.34 ± 0.17	0.21 ± 0.15	0.25 ± 0.18	0.3 ± 0.22	0.39 ± 0.22	0.19 ± 0.14	0.23 ± 0.15	0.3 ± 0.17	0.36 ± 0.16
AvgProb	0.17 ± 0.12	0.21 ± 0.14	0.26 ± 0.17	0.38 ± 0.15	0.17 ± 0.13	0.22 ± 0.15	0.26 ± 0.19	0.37 ± 0.18	0.16 ± 0.13	0.22 ± 0.15	0.29 ± 0.16	0.37 ± 0.16
pS-SDM+PRR	0.15 ± 0.16	0.17 ± 0.18	0.26 ± 0.23	0.28 ± 0.25	0.14 ± 0.17	0.18 ± 0.2	0.22 ± 0.23	0.29 ± 0.29	0.14 ± 0.16	0.17 ± 0.19	0.26 ± 0.21	0.28 ± 0.22
MEM+PRR	0.14 ± 0.16	0.17 ± 0.21	0.24 ± 0.22	0.32 ± 0.26	0.14 ± 0.17	0.18 ± 0.2	0.22 ± 0.23	0.34 ± 0.29	0.15 ± 0.16	0.19 ± 0.19	0.25 ± 0.21	0.32 ± 0.22

8

9 C) Percentage correct classified (PCC)

Thresholding Approach	SSCV				ID				CCV			
	FO	GL	BF	GH	FO	GL	BF	GH	FO	GL	BF	GH
Max.Kappa	0.91 ± 0.09	0.9 ± 0.07	0.78 ± 0.1	0.83 ± 0.08	0.9 ± 0.09	0.89 ± 0.07	0.77 ± 0.09	0.83 ± 0.1	0.88 ± 0.09	0.87 ± 0.06	0.76 ± 0.08	0.82 ± 0.07
MCE05	0.85 ± 0.15	0.79 ± 0.07	0.66 ± 0.11	0.82 ± 0.08	0.88 ± 0.07	0.87 ± 0.05	0.77 ± 0.09	0.82 ± 0.08	0.59 ± 0.14	0.68 ± 0.09	0.64 ± 0.1	0.73 ± 0.08
Max.TSS	0.85 ± 0.08	0.84 ± 0.07	0.73 ± 0.11	0.83 ± 0.09	0.86 ± 0.1	0.85 ± 0.08	0.73 ± 0.11	0.81 ± 0.09	0.77 ± 0.07	0.79 ± 0.05	0.72 ± 0.08	0.79 ± 0.06
SensSpec	0.79 ± 0.1	0.84 ± 0.04	0.73 ± 0.09	0.81 ± 0.08	0.87 ± 0.09	0.86 ± 0.07	0.73 ± 0.1	0.81 ± 0.08	0.79 ± 0.07	0.81 ± 0.05	0.73 ± 0.07	0.8 ± 0.06
Opt.ROC	0.86 ± 0.08	0.86 ± 0.06	0.72 ± 0.11	0.82 ± 0.05	0.87 ± 0.1	0.86 ± 0.07	0.74 ± 0.1	0.81 ± 0.09	0.79 ± 0.07	0.81 ± 0.05	0.73 ± 0.08	0.8 ± 0.06
Preval	0.92 ± 0.08	0.91 ± 0.05	0.79 ± 0.08	0.83 ± 0.06	0.9 ± 0.08	0.89 ± 0.06	0.77 ± 0.1	0.83 ± 0.09	0.9 ± 0.08	0.89 ± 0.06	0.77 ± 0.08	0.82 ± 0.07
AvgProb	0.71 ± 0.07	0.67 ± 0.08	0.64 ± 0.11	0.74 ± 0.08	0.69 ± 0.08	0.69 ± 0.07	0.64 ± 0.11	0.73 ± 0.08	0.69 ± 0.07	0.69 ± 0.07	0.66 ± 0.09	0.73 ± 0.07
pS-SDM+PRR	0.93 ± 0.09	0.88 ± 0.07	0.76 ± 0.08	0.86 ± 0.09	0.91 ± 0.1	0.89 ± 0.08	0.77 ± 0.11	0.83 ± 0.1	0.91 ± 0.1	0.89 ± 0.08	0.78 ± 0.09	0.84 ± 0.1
MEM+PRR	0.92 ± 0.1	0.9 ± 0.09	0.8 ± 0.09	0.84 ± 0.09	0.91 ± 0.1	0.89 ± 0.09	0.79 ± 0.1	0.86 ± 0.08	0.91 ± 0.1	0.89 ± 0.08	0.79 ± 0.08	0.86 ± 0.09

10

11

12 D) Sensitivity

Thresholding Approach	SSCV				ID				CCV			
	FO	GL	BF	GH	FO	GL	BF	GH	FO	GL	BF	GH
Max.Kappa	0.31 ± 0.18	0.36 ± 0.24	0.47 ± 0.21	0.51 ± 0.26	0.27 ± 0.21	0.32 ± 0.25	0.46 ± 0.24	0.56 ± 0.26	0.32 ± 0.18	0.37 ± 0.18	0.5 ± 0.2	0.52 ± 0.22
MCE05	0.45 ± 0.13	0.61 ± 0.11	0.66 ± 0.18	0.73 ± 0.12	0.35 ± 0.17	0.42 ± 0.17	0.42 ± 0.23	0.57 ± 0.19	0.68 ± 0.05	0.67 ± 0.07	0.7 ± 0.09	0.71 ± 0.11
Max.TSS	0.44 ± 0.28	0.55 ± 0.15	0.62 ± 0.14	0.67 ± 0.15	0.46 ± 0.28	0.52 ± 0.27	0.59 ± 0.22	0.66 ± 0.25	0.57 ± 0.11	0.59 ± 0.11	0.61 ± 0.14	0.63 ± 0.15
SensSpec	0.55 ± 0.25	0.49 ± 0.15	0.58 ± 0.18	0.65 ± 0.13	0.42 ± 0.26	0.49 ± 0.25	0.58 ± 0.22	0.67 ± 0.25	0.54 ± 0.1	0.55 ± 0.11	0.58 ± 0.14	0.63 ± 0.12
Opt.ROC	0.44 ± 0.2	0.54 ± 0.14	0.57 ± 0.21	0.67 ± 0.14	0.44 ± 0.26	0.49 ± 0.26	0.58 ± 0.22	0.65 ± 0.24	0.55 ± 0.1	0.56 ± 0.11	0.59 ± 0.14	0.62 ± 0.14
Preval	0.28 ± 0.17	0.28 ± 0.21	0.46 ± 0.21	0.55 ± 0.25	0.26 ± 0.2	0.32 ± 0.24	0.49 ± 0.22	0.53 ± 0.25	0.24 ± 0.18	0.29 ± 0.18	0.47 ± 0.19	0.5 ± 0.23
AvgProb	0.76 ± 0.11	0.82 ± 0.1	0.77 ± 0.12	0.83 ± 0.1	0.79 ± 0.11	0.85 ± 0.1	0.78 ± 0.17	0.85 ± 0.13	0.76 ± 0.1	0.81 ± 0.06	0.78 ± 0.1	0.84 ± 0.07
pS-SDM+PRR	0.21 ± 0.28	0.26 ± 0.32	0.43 ± 0.36	0.42 ± 0.39	0.21 ± 0.28	0.28 ± 0.33	0.44 ± 0.36	0.45 ± 0.4	0.21 ± 0.27	0.26 ± 0.3	0.45 ± 0.33	0.41 ± 0.36
MEM+PRR	0.21 ± 0.28	0.3 ± 0.31	0.38 ± 0.36	0.41 ± 0.38	0.21 ± 0.29	0.28 ± 0.33	0.38 ± 0.35	0.43 ± 0.38	0.21 ± 0.27	0.28 ± 0.3	0.39 ± 0.32	0.42 ± 0.34

13

14

15

16 E) Specificity

Thresholding Approach	SSCV				ID				CCV			
	FO	GL	BF	GH	FO	GL	BF	GH	FO	GL	BF	GH
Max.Kappa	0.95 ± 0.09	0.91 ± 0.07	0.82 ± 0.12	0.87 ± 0.1	0.94 ± 0.07	0.93 ± 0.07	0.84 ± 0.12	0.88 ± 0.1	0.91 ± 0.09	0.93 ± 0.07	0.81 ± 0.12	0.85 ± 0.1
MCE05	0.87 ± 0.15	0.81 ± 0.04	0.78 ± 0.11	0.79 ± 0.09	0.93 ± 0.03	0.91 ± 0.03	0.86 ± 0.1	0.88 ± 0.08	0.57 ± 0.15	0.91 ± 0.03	0.61 ± 0.13	0.7 ± 0.1
Max.TSS	0.88 ± 0.09	0.79 ± 0.07	0.73 ± 0.1	0.83 ± 0.09	0.87 ± 0.1	0.86 ± 0.09	0.75 ± 0.13	0.84 ± 0.09	0.78 ± 0.08	0.86 ± 0.09	0.74 ± 0.1	0.8 ± 0.08
SensSpec	0.89 ± 0.07	0.82 ± 0.06	0.77 ± 0.1	0.81 ± 0.09	0.89 ± 0.09	0.88 ± 0.07	0.75 ± 0.12	0.84 ± 0.09	0.81 ± 0.07	0.88 ± 0.07	0.76 ± 0.09	0.82 ± 0.07
Opt.ROC	0.84 ± 0.08	0.87 ± 0.08	0.74 ± 0.11	0.82 ± 0.09	0.88 ± 0.09	0.88 ± 0.08	0.76 ± 0.13	0.84 ± 0.09	0.8 ± 0.07	0.88 ± 0.08	0.76 ± 0.09	0.82 ± 0.07
Preval	0.92 ± 0.07	0.95 ± 0.05	0.85 ± 0.13	0.85 ± 0.09	0.94 ± 0.08	0.94 ± 0.06	0.82 ± 0.13	0.87 ± 0.11	0.94 ± 0.08	0.94 ± 0.06	0.84 ± 0.11	0.87 ± 0.09
AvgProb	0.69 ± 0.08	0.68 ± 0.06	0.63 ± 0.12	0.72 ± 0.09	0.68 ± 0.08	0.68 ± 0.07	0.6 ± 0.12	0.71 ± 0.09	0.68 ± 0.08	0.68 ± 0.07	0.62 ± 0.09	0.69 ± 0.08
pS-SDM+PRR	0.92 ± 0.15	0.91 ± 0.12	0.81 ± 0.19	0.86 ± 0.19	0.93 ± 0.15	0.92 ± 0.12	0.8 ± 0.23	0.84 ± 0.21	0.93 ± 0.14	0.92 ± 0.12	0.82 ± 0.19	0.86 ± 0.18
MEM+PRR	0.93 ± 0.14	0.93 ± 0.1	0.85 ± 0.18	0.91 ± 0.15	0.93 ± 0.15	0.92 ± 0.13	0.85 ± 0.2	0.89 ± 0.16	0.93 ± 0.14	0.92 ± 0.13	0.86 ± 0.17	0.88 ± 0.15

17

18

19

20