

## RESEARCH ARTICLE



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

Daniel Scherrer<sup>1</sup> | Manuela D'Amen<sup>1</sup> | Rui F. Fernandes<sup>1</sup> |  
Rubén G. Mateo<sup>1,2</sup> | Antoine Guisan<sup>1,3</sup>

<sup>1</sup>Department of Ecology and Evolution, University of Lausanne, Biophore, Lausanne, Switzerland

<sup>2</sup>ETSI de Montes, Forestal y del Medio Natural, Universidad Politécnica de Madrid, Madrid, Spain

<sup>3</sup>Institute of Earth Surface Dynamics, University of Lausanne, Géopolis, Lausanne, Switzerland

## Correspondence

Daniel Scherrer, Department of Ecology and Evolution, University of Lausanne, CH-1015 Lausanne, Switzerland.  
Email: daniel.scherrer@unil.ch

## Funding information

H2020 Marie Skłodowska-Curie Actions; Marie Curie Intra-European Fellowship within the 7th European Community Framework Programme, Grant/Award Number: ACONITE, PIEF-GA-2013-622620; Schweizerischer Nationalfonds zur Förderung der Wissenschaftlichen Forschung, Grant/Award Number: 31003A-1528661

Handling Editor: Jana McPherson

## Abstract

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

## KEYWORDS

community cross-validation, community modeling, ecological niche models, species distribution modeling, species distribution models, thresholding methods

## 1 | INTRODUCTION

Past and future environmental changes may not only lead to shifts in species distributions (e.g. Dullinger et al., 2012; Parmesan & Yohe, 2003; Thuiller, Lavorel, Araújo, Sykes, & Prentice, 2005), but also to changes in species assemblages and interactions (e.g. Alexander, Diez, Hart, & Levine, 2016; Blois, Zarnetske, Fitzpatrick, & Finnegan, 2013; Nogues-Bravo & Rahbek, 2011; Van der Putten, Macel, & Visser, 2010). Information about communities, here defined as a taxonomic assemblage of distinct populations of species that co-occur in a given space at a given time (Begon, Harper, & Townsend, 1996), is therefore essential to make informed decisions for conservation prioritisation (D'Amen et al., 2011; Guisan et al., 2013; Mateo, de la Estrella, Felicísimo, Munoz, & Guisan, 2013) and to create biodiversity indices (e.g. Essential Biodiversity Variables; Pereira et al., 2013) for policy decisions (Fleishman, Noss, & Noon, 2006; Granger et al., 2015).

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

While there is a vast and now long-standing literature on advances and limitations of single species predictions (e.g. Elith & Leathwick, 2009; Guisan & Thuiller, 2005; Guisan et al., 2006; Maggini, Lehmann, Zimmermann, & Guisan, 2006; Meier et al., 2010; Merow et al., 2014; Zimmermann, Edwards, Graham, Pearman, & Svenning, 2010), studies exploring how to improve community predictions based on aggregated information from individual SDMs emerged more recently (e.g. Benito, Cayuela, & Albuquerque, 2013; Cord, Klein, Gernandt, la Rosa, & Dech, 2014; Mateo, Felicísimo, Pottier, Guisan, & Munoz, 2012; Mod, le Roux, Guisan, & Luoto, 2015; but see Ferrier, Drielsma, Manion, & Watson, 2002). A fundamental difference among the proposed solutions is whether to maintain the information on species composition in the final predictions. For instance, the simple sum of probabilities of individual SDM predictions usually gives better estimates of species richness, but the information on species identity is lost (Calabrese, Certain, Kraan,

& Dormann, 2014; Dubuis et al., 2011). Therefore, predictions of community composition have mainly been achieved so far by thresholding the individual continuous SDM predictions (e.g. probability or suitability index) to obtain binary maps (Liu, White, & Newell, 2013) and then stacking the latter at the assemblage level (e.g. D'Amen, Dubuis, et al., 2015; D'Amen, Pradervand, & Guisan, 2015; Pottier et al., 2013).

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

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

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

Surprisingly, studies aiming to test and improve S-SDM have used very different approaches to evaluate the predicted assemblages (Cord et al., 2014; Hespanhol, Cezón, Felicísimo, Muñoz, & Mateo, 2015; Pouteau et al., 2015; Thuiller, Pollock, Gueguen, & Münkemüller, 2015; Zurell, Zimmermann, Sattler, Nobis, & Schröder, 2016) and this evaluation aspect of the community modelling procedure has not yet received all the attention it deserves.

In most studies, assemblage predictions are not adequately evaluated because the data used for the evaluation were already used for individual model fitting, not allowing anymore a correct cross-validation at the community level. Ideally, the best evaluation method should use spatial or temporal independent data (Elith et al., 2006; Guisan et al., 2017), but if not available, an appropriate cross-validation approach should at least be set up.

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

## 2 | MATERIALS AND METHODS

### 2.1 | Community data and environmental variables

#### 2.1.1 | Study area

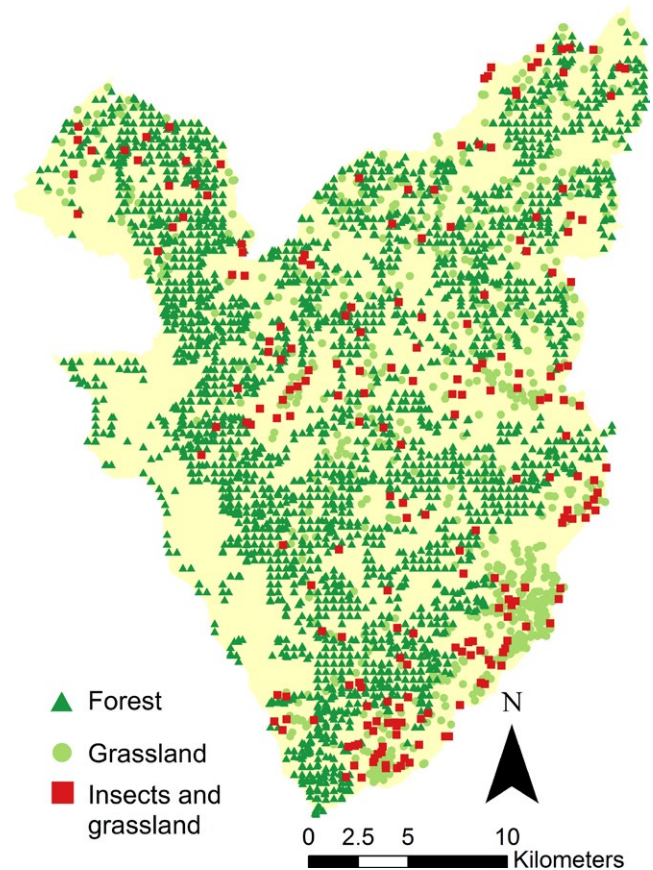
The data on all taxa were collected within the same study area located in the western Swiss Alps of the canton Vaud (Figure 1; 46°10' to 46°30' N; 6°50' to 7°10' E), covering an area of ca. 700 km<sup>2</sup>, with elevation ranging from 375 to 3,210 m a.s.l. and forested areas up to 1,900 m a.s.l. For centuries, agriculture (farming and pasturing) has maintained grasslands among forests and altered the position of the treeline. The highly variable topography and diverse land use of the study area, in combination with our high-resolution environmental data (25 × 25 m cell size), provide a huge range of complex species-environment relationships to test our modelling framework.

#### 2.1.2 | Plant data

The forest data were part of a forest inventory of the canton Vaud conducted between 1988 and 2002 (mostly 1990 to 1994) and consisted of 3,076 sites. The forest sites were distributed on a 400 m grid all across the forested area of the canton and had a circular area of 314 m<sup>2</sup> (Figure 1; for details see Hartmann, Fouvry, & Horisberger, 2009). In total, 703 plant species were recorded, but only 312 (44%) had enough occurrence data (>20 occurrences) across the dataset for modelling purposes (see Table 1 for more detailed statistics on the datasets).

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

To predict the distribution of the plant species, we used five environmental variables: growing degree-day (above 0°C), moisture



**FIGURE 1** Map of the study area with the forested sites (dark green triangles,  $N = 3,076$ ), the grassland sites (light green circles and red squares,  $N = 903$ ) and the insect sites (red squares, butterflies  $N = 192$ , grasshoppers  $N = 202$ )

index over the growing season (difference between precipitation and potential evapotranspiration), the sum of potential solar radiation over the year, slope (in degrees), and topographic position (unitless, indicating the ridges and valleys). All these variables were at a 25 m resolution and have been shown to be useful predictors for plant species in mountain environments (see D'Amen, Dubuis, et al., 2015; Dubuis et al., 2011; Scherrer, Massy, Meier, Vittoz, & Guisan, 2017 for details on predictors).

#### 2.1.3 | Insect data

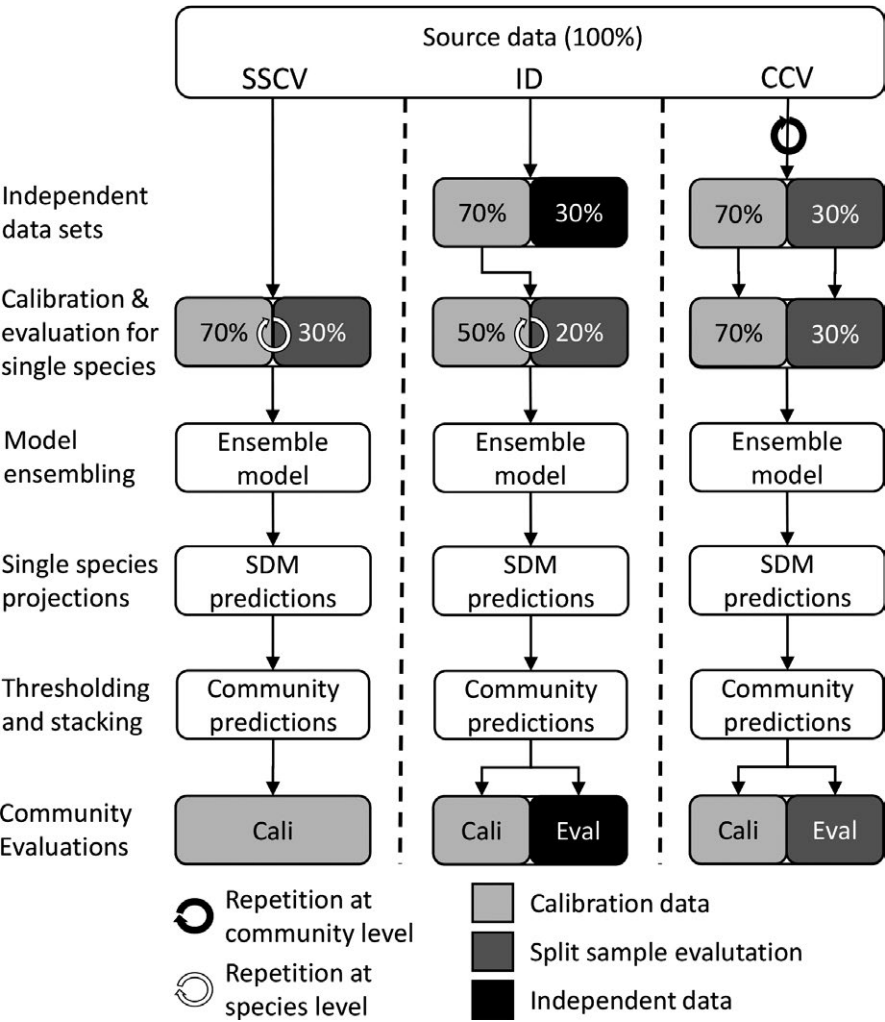
Data on butterflies and grasshoppers were, respectively, collected in 192 and 202 squares of 50 × 50 m across all the elevational range of the study area (Figure 1; see Pellissier et al., 2012; Pradervand et al., 2013; for more information). In total, 131 butterfly and 41 grasshopper species were observed,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

For our SDMs, we used the same predictors as D'Amen, Pradervand, et al. (2015): four bioclimatic variables (solar radiation, summer temperature, annual degree-days and annual average number of frost days during the growing season), an index of vegetation

**TABLE 1** Basic statistics of the datasets used for the case study and the evaluation metrics (AUC) for the individual species distribution models, using three different community evaluation approaches

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

CCV, community cross-validation; ID, independent data; SSCV, single species cross-validation.



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

productivity, that is, normalised difference vegetation index (as proxies for trophic resources), and the distance to forest. These variables were selected as they are not highly correlated (<0.7; Dormann et al., 2013) and considered ecologically important for insects (e.g. Hawkins & Porter, 2003; Turner, Gatehouse, & Corey, 1987).

## 2.2 | The modelling framework

Our modelling framework used three different S-SDM-based community modelling pathways (“single species cross-validation”,

“independent data” and “community cross-validation”) representing the most commonly reported practices in the literature (see Figure 2 and “Evaluating community predictions” section).

### 2.2.1 | Single species modelling, thresholding and evaluation

Individual species models were run by generalised linear models (GLM; McCullagh & Nelder, 1989), generalised additive models (GAM; Hastie & Tibshirani, 1990), random forest (RF; Breiman,

2001) and boosted regression trees (BRT; Elith, Leathwick, & Hastie, 2008). Models for species with more than 50 occurrences were fitted by simple SDMs using all five selected predictors, followed by a weighted (AUC) ensemble forecast (Marmion, Parviainen, Luoto, Heikkinen, & Thuiller, 2009). Species having only between 20 and 50 occurrence records were fitted by an ensemble bivariate approach optimised for rare or under-sampled species (Breiner, Guisan, Bergamini, & Nobis, 2015; Lomba et al., 2010): individual models were calibrated on bivariate combinations of the selected predictors with all four modelling techniques, followed by a consensus forecast from all the resulting “small models” weighted by their AUC scores. We used a repeated split-sample procedure ( $N = 25$ ) for model evaluation, followed by a weighted (AUC) ensemble forecast (across techniques and split-sample runs).

The projected probability outputs of the ensemble models were binarised using two thresholding schemes: (a) *species-specific-thresholds* (a single threshold calculated for each species) and (b) *site-specific-thresholds* (differing for each site on the basis of additional community information, that is, species richness predictions). We selected seven different species-specific thresholding techniques, which can be classified into four major groups: single-index based, sensitivity and specificity combined, model-building data-only-based and predicted probability-based (see Supporting Information Table S1; Liu et al., 2005; Nenzen & Araujo, 2011 for details on classification). As the thresholding techniques showed minimal within-group variance (see Supporting Information Figures S1 and S2), we decided to only present the results for one thresholding technique per group in the main manuscript. The chosen techniques were as follows: Cohen's Kappa maximisation approach (*Max.Kappa*; single-index based), TSS maximisation approach (*Max.TSS*, sensitivity and specificity combined), observed prevalence (*Obs.Preval*; model-building data-only-based approach), and average probability approach (*AvgProb*; predicted probability-based approach; for details on techniques see Supporting Information Table S1). In addition, we applied two site-thresholds (community-based approaches), using species richness (SR) predictions in combination with a probability ranking rule (PRR). These methods selected a number of species equal to the predicted SR on the basis of decreasing probabilities of presence calculated by the SDMs (D'Amen, Dubuis, et al., 2015; D'Amen, Pradervand, et al., 2015). Therefore, the species with the highest probabilities in a site are selected (considered present) in decreasing order until the SR predicted for the site is reached. The SR predictions were derived by either summing the per site probabilities of individual SDMs, obtaining a prediction of richness for each site (pS-SDM; Dubuis et al., 2011) or by a macroecological model (MEM; see D'Amen, Pradervand, et al., 2015 for details), directly modelling the richness of the sites. As results from the two site thresholds were concordant, we only show here the former (pS-SDM + PRR).

To evaluate the threshold independent performance of our individual species models, the area under the curve of a

Receiver-Operating Characteristic (ROC) plot (AUC; Fielding & Bell, 1997) was calculated based on a repeated split sampling cross-validation (Thuiller, Georges, & Engler, 2013). Additionally, based on our independent/cross-validation data, we calculated five threshold-dependent metrics for each thresholding technique: the overall accuracy (PCC; i.e. proportion of correctly classified presence and absences; Fielding & Bell, 1997), sensitivity (proportion of correctly predicted presences), specificity (proportion of correctly predicted absences), the true skill statistic (i.e. [(sensitivity + specificity) - 1]; TSS; Allouche et al., 2006) and Cohen's Kappa (Kappa; i.e. overall accuracy but corrected for chance performance; Cohen, 1968).

## 2.2.2 | Evaluating community predictions

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

- The “single species cross-validation” (SSCV) approach (Figure 2) has not fully “unused/independent” data for community evaluation (i.e. sites not used for the calibration of any single species). Here, in the process of the cross-validation of all individual SDMs (i.e. across all species), different sites are selected at each resampling iteration and for each species, so that all sites are most likely used in at least one split-sampling run and their information incorporated in the final ensemble model. This approach cannot thus be considered based on fully independent data. The SSCV approach has been to date the most common way to model and evaluate communities predictions based on S-SDMs (Figure 2; for example, Dubuis et al., 2011; Calabrese et al., 2014; D'Amen, Pradervand, et al., 2015; Distler et al., 2015). As no independent data is set aside for community evaluation, this approach usually gets evaluated with all the sites used for calibration. However, to avoid bias in the results due to different numbers of evaluation sites, we evaluated the SSCV approach only on 30% of the available sites (identical to the ID and CCV approach below).
- The (spatial or temporal) “independent data” (ID) approach (Figure 2) starts with two completely independent datasets. One is used for the calibration of the SDMs (i.e. 70% of the sites) and the other set is used (only) to evaluate the performance of the community predictions (i.e. 30% of the sites; Figure 2; for example, Benito et al., 2013; Cord et al., 2014; D'Amen, Dubuis, et al., 2015; Pottier et al., 2013; Zurell et al., 2016).
- The “community cross-validation” (CCV) approach (Figure 2) uses a repeated split sampling of sites (100 repetitions) dividing the available sites into calibration (70%) and evaluation sets (30%) to perform all the modelling procedure from the single species prediction to the community assembly (Figure 2). In



**TABLE 2** Pearson Correlation of single species and community evaluation statistics. Correlations of the single species evaluation metrics and their corresponding community evaluation metric are highlighted in bold

Single species	Community metrics					Sørensen similarity	SR deviation
	Accuracy	Sensitivity	Specificity	KAPPA	TSS		
Accuracy	<b>1.00***</b>	−0.37*	0.95***	0.70***	0.37*	0.37*	−0.58***
Sensitivity	−0.36**	<b>0.93***</b>	−0.54***	0.01 n.s.	0.56***	0.18 n.s.	−0.44***
Specificity	0.97***	−0.53***	<b>0.99***</b>	0.64***	0.20 n.s.	0.31*	−0.63***
KAPPA	0.41**	0.50*	0.27 *	<b>0.79 ***</b>	0.72***	0.82***	−0.3*
TSS	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, \* $p < 0.05$ , \*\* $p < 0.01$ , \*\*\* $p < 0.001$ )

contrast to the previous ID pathway (above), which only uses one (spatial or temporal) fixed independent evaluation dataset, in the CCV approach all SDMs are fitted at each split-sample iteration using the same training and test sets for all species, thus minimising the risk of bias in the evaluation data (i.e. if the training and test sets differ across species, as in the ID approach). This repeated cross-validation also allows the estimation/simulation of confidence intervals for community predictions instead of just a single value per community. To our knowledge, no study used this community cross-validation method so far.

To compare the community model performance among thresholding techniques and modelling pathways, we calculated eight different community agreement metrics: (a) the deviation of the predicted from the observed species richness (SR.deviation), (b) the proportion of species correctly predicted as present (community sensitivity), (c) the proportion of species correctly predicted as absent (community specificity), (d) community accuracy (PCC; i.e. the percent correctly classified species, present or absent), (e) the community TSS (here measured for a site across all species, rather than for a species across all sites as in single SDM evaluation; Pottier et al., 2013), (f) the community kappa (same as for TSS, for a site across species; Pottier et al., 2013), and (g) the Sørensen similarity (Sørensen, 1948).

### 2.2.3 | Correlation of single species and community evaluation metrics

For each combination of dataset, the modelling pathway and thresholding method, ( $4 \times 3 \times 9 = 108$ ) we calculated the average evaluation metric for all five single species metrics and all seven community metrics. We then calculated the Spearman correlation of all possible combinations of our five single species and seven community evaluation metrics. The resulting correlation matrix tells us if methods (modelling pathways or thresholding methods) that yield the highest scores in a certain single species metric also yield the highest score in the corresponding community evaluation metric.

## 3 | RESULTS

### 3.1 | Performance of individual SDMs

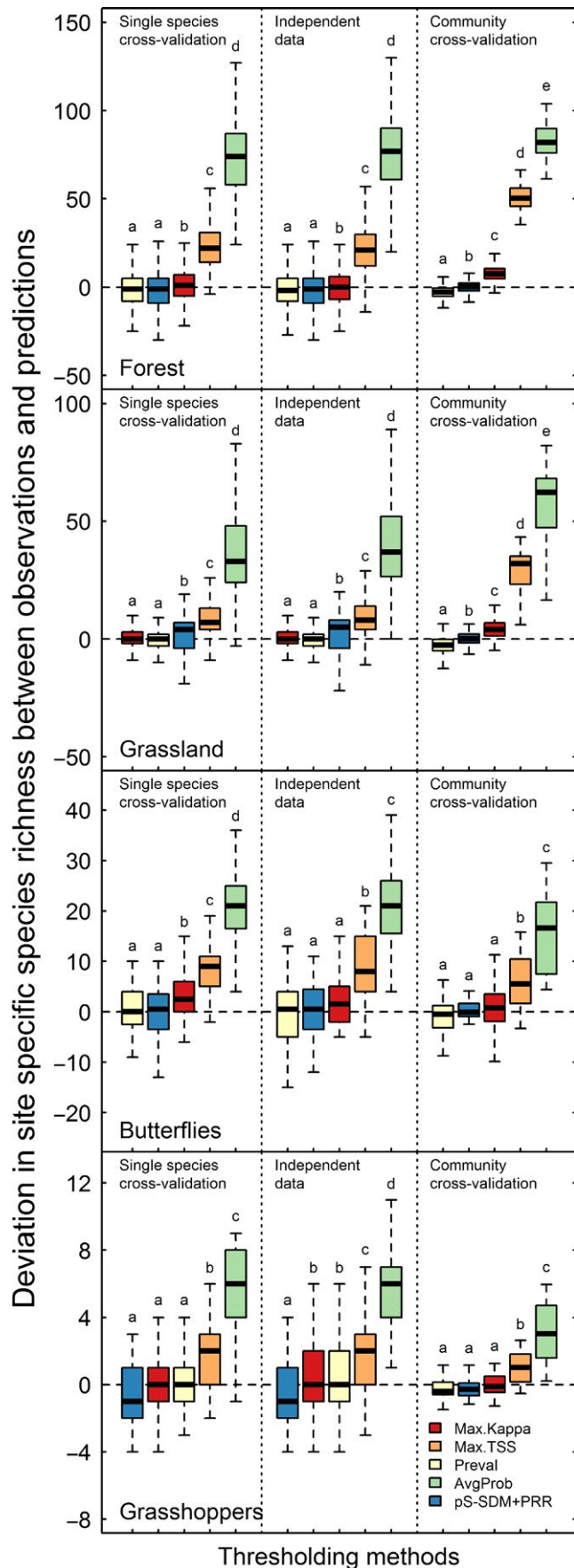
As expected, the evaluation scores of the individual SDMs were similar to earlier studies published with the same data (D'Amen, Dubuis, et al., 2015; D'Amen, Pradervand, et al., 2015; Scherrer et al., 2017) and their performance was not affected by the chosen community evaluation approach (Table 1, Supporting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

### 3.2 | Correlation of single species and community evaluation metrics

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

### 3.3 | Species richness and compositional similarity

The deviation in species richness between observed and predicted communities was strongly dependent on the chosen thresholding method (Figure 3). The thresholding approach that uses the average predicted probability (*AvgProb*) showed the highest amount of over-prediction followed by the combined sensitivity and specificity



**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 25th to the 75th percentile and the whiskers indicate  $\pm 2$  SD. Letters above the boxplots indicate significant differences (Wilcoxon rank sum test,  $p < 0.05$ )

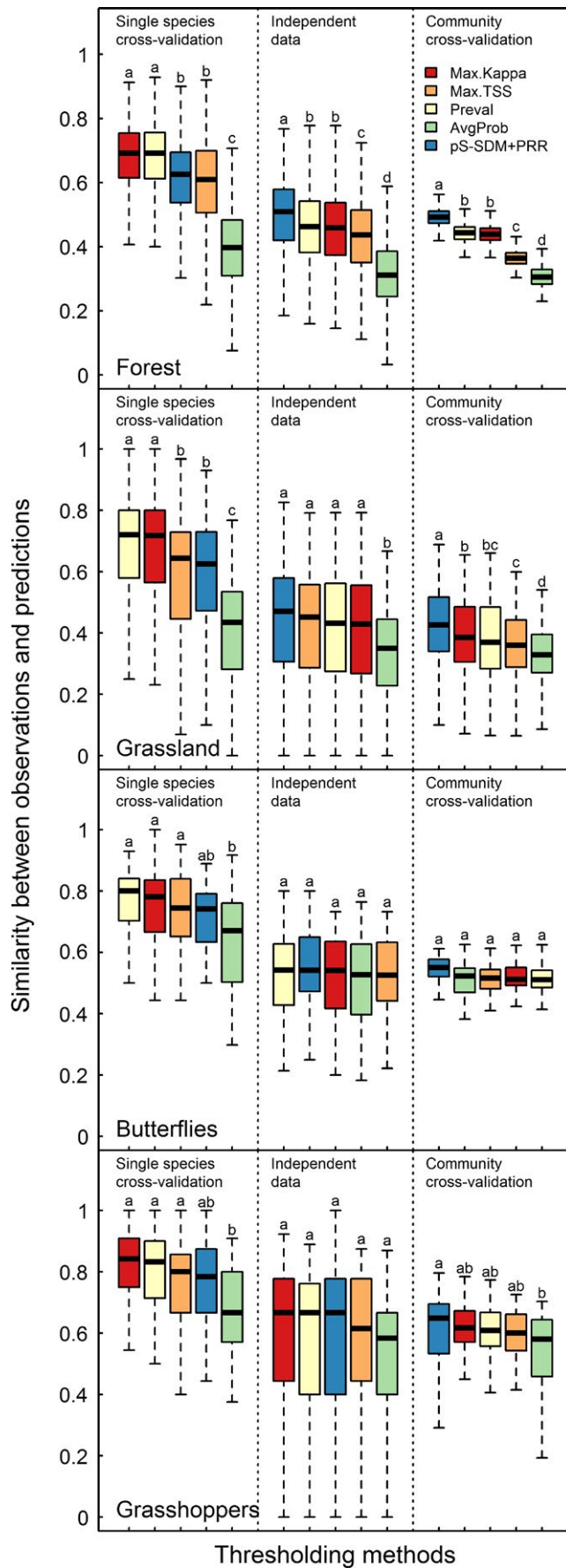
approach (*Max.TSS*). The other three thresholding methods (*Obs.Preval*, *Max.Kappa* and *pS-SDM+PRR*) performed in a similar way and showed overall no tendency to over-predict species richness. There were no significant differences between the three modelling pathways for any of the studied taxa (Figure 3). The absolute number of over-predicted species was strongly related to the average number of species per plot (SR) and therefore differed among the taxa (Figure 3). However, when corrected for the differences in SR the over-prediction did not significantly vary anymore across taxa.

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

## 4 | DISCUSSION

### 4.1 | Do the most established thresholds for single species work as well for community predictions?

In this paper, we asked if the most established methods for single species thresholding are also the optimal choice for making predictions at the community level and if there is a direct link between the individual species predictions and the corresponding community metrics. Our results confirm the existence of such a link for



**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 25th to the 75th percentile and the whiskers indicate  $\pm 2$  SD. Letters above the boxplots indicate significant differences (Wilcoxon rank sum test,  $p < 0.05$ )

single-index based metrics such as sensitivity, specificity and accuracy. However, these results should be interpreted with caution as maximising sensitivity or specificity can simply be achieved by predicting the species as present or absent (respectively) everywhere. In our study system, most of the modelled species have a low prevalence (i.e. are absent at most sites), thus accuracy (PCC) can often be improved by predicting the species as “absent” nearly everywhere.

The two most commonly used community evaluation metrics, Sørensen similarity index and deviation in species richness, were only weakly correlated with most evaluation metrics used for individual species. The most established thresholding methods for individual species predictions (i.e. *Max.TSS*, *Opt.ROC*, *SenSpec*) did show lower performance when applied to community-level predictions. This is likely due to the fact that both TSS and ROC try to find the best trade-off between sensitivity and specificity (Guisan et al., 2017). As most of the species have a prevalence far below 50% (i.e. are absent in many more sites than present), adding a few more presences might have a big effect on the sensitivity (by increasing the chance of finding the few real presences) but only marginally affects the specificity. By definition, increasing sensitivity also increases TSS, but with the drawback of a slight over-prediction. While this might not matter much on a single species basis, for community-level predictions, the over-prediction will accumulate when summing binarised maps across all species, leading to the often observed over-estimation of species richness in S-SDMs (e.g. Dubuis et al., 2011; Mateo et al., 2012; Pineda & Lobo, 2009; Pottier et al., 2013; Pouteau et al., 2015; Zurell et al., 2016). It is important to remark, that in the rare case of an ecosystem mostly comprising of widespread species (i.e. prevalence  $>50\%$ ) this will turn into the opposite as TSS and ROC will optimise absences leading to an underestimation of species richness. The strength of the over/under prediction bias is therefore linked to the prevalence distribution of the modelled species assemblages. However, in the vast majority of natural systems, both the site SR and the regional species pool are driven by a large number of rare (low prevalence species) compared to a few widespread species (Magurran & Henderson, 2003; Preston, 1948).

The community-based thresholding methods based on the selection of the most probable species (through a probability ranking) up to the predicted site richness (*MEM+PRR*, *pS-SDM+PRR*) can overcome this problem, because they are able to constrain species predictions based on a different value of species richness in each site (i.e. making them site-specific thresholding methods). Therefore, these methods prevent over-prediction while still allowing the analyses of species composition. Our results thus support the conclusion



that, when the final goal is to optimise community composition, community-thresholding methods are the best option. Yet, as discussed in the next section, two single-species thresholding methods – *maximised Kappa* and *observed prevalence* – also showed good results for predicting communities (close to the community-based approaches). However, as community-based thresholds combine the optimisation of species richness prediction and a probability ranking rule (PRR), they would always select the species with the highest predicted probabilities in each site (D'Amen, Pradervand, et al., 2015). This could seem logic and straightforward, but there might be a bias when the species in the community have varying prevalence (D'Amen, Mateo, et al., 2017). In fact, the maximum predicted probability is depending on the prevalence of the species, thus the common species will tend to always have greater maximum predicted probabilities than rare species and, as a result, will be considered present an over-proportionate number of time in the final community compositions. This bias will produce high similarity scores (Sørensen index) in the prediction evaluation, as the most common species are correctly predicted in most sites. However, the drawback is that the rarest species will be often omitted in the community predictions, which can be, for instance, problematic if the final goal of the modelling exercise has conservation implications.

## 4.2 | Is there a “best” threshold for community S-SDMs?

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

A similar ranking of thresholding methods was overall observed across taxonomic group within a given modelling pathway, while among the pathways there were clear shifts in the ranking of thresholding methods: with no independent community evaluation data (SSCV), the *Obs.Preval* and *Max.Kappa* threshold showed superior results, while the pathways using independent community evaluation data (ID and CCV) indicated the community-based thresholding to be superior (*pS-SDM+PRR*). This observation is in line with published literature, where studies not using independent community data usually report a good performance of single species optimisations methods (e.g. D'Amen, Pradervand, et al., 2015; Distler et al., 2015;

Thuiller et al., 2015), while studies using independent data usually have better results using community constraints (e.g. D'Amen, Dubuis, et al., 2015). Yet, it is remarkable to notice that, although previously much criticised in the literature (e.g. Allouche et al., 2006; McPherson, Jetz, & Rogers, 2004), maximised Kappa (together here with the observed prevalence) did indeed perform well as a thresholding method for predicting both single species and communities, being nearly always superior to the sensitivity-specificity thresholding methods supporting earlier findings of Manel, Williams, and Ormerod (2001).

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

## 4.3 | Summing up: How to evaluate community predictions correctly?

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

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

The third approach (CVV), newly presented here, repeats the ID approach a large number of times within a cross-validation procedure at the community level (no example of this approach known in the literature). By doing this, the risk of bias in the evaluation data, inherent to the selection of a single evaluation dataset, is minimised compared to the simple ID approach. Additionally, the repeated cross-validation allows the assessments of uncertainty and confidence intervals around the community predictions' performance

metrics. However, as this approach selects the same sites for all species, its application is only possible under specific circumstances. First, all the species data need to be collected in the same sites (i.e. true “community data”). Second, as this approach leads to an unequal number of presences/absences between different cross-validation runs for the same species, it can lead to models failing for very rare (low sample size) species in some of the cross-validation runs if not enough presence sites are selected in the training set.

According to our results and despite the potential limitations, we advise the use of the proposed community cross-validation approach (CCV) to evaluate community models in future studies. In fact, we clearly showed that the common practice of evaluating the community predictions on the same dataset used for calibration process (SSCV) leads to overoptimistic estimations of model performance. In the commonest case of unavailability of truly spatial (i.e. different region) or temporal (i.e. different sampling period) independent data, often independent datasets are “created” by randomly splitting the initial dataset in two parts. However, we advocate against this practise and instead promote the community cross-validation approach, which minimises the artefacts of randomly splitting the initial data and allows the estimation of uncertainty associated with the community evaluation metrics.

## ACKNOWLEDGEMENTS

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

## AUTHORS' CONTRIBUTIONS

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

## DATA ACCESSIBILITY

A generalised version of the community cross-validation algorithm is available in the *ecospat* R package (Cola et al., 2016) on GitHub (*ecospat.CCV*; <https://doi.org/10.5281/zenodo.1287805>). All species and environmental data are available on Dryad Digital Repository <https://doi.org/10.5061/dryad.28d4k> (Grassland species and environmental predictors for plants; Guisan, Dubuis, & Vittoz, 2011) and <https://doi.org/10.5061/dryad.nf925> (forest, insect species and environmental predictors for insects; Guisan et al., 2018).

## ORCID

Daniel Scherrer  <http://orcid.org/0000-0001-9983-7510>

## REFERENCES

- Alexander, J. M., Diez, J. M., Hart, S. P., & Levine, J. M. (2016). When climate reshuffles competitors: A call for experimental macroecology. *Trends in Ecology & Evolution*, 31, 831–841. <https://doi.org/10.1016/j.tree.2016.08.003>
- Allouche, O., Tsoar, A., & Kadmon, R. (2006). Assessing the accuracy of species distribution models: Prevalence, kappa and the true skill statistic (TSS). *Journal of Applied Ecology*, 43, 1223–1232. <https://doi.org/10.1111/j.1365-2664.2006.01214.x>
- Begon, M., Harper, J. L., & Townsend, C. R. (1996). *Ecology: Individuals populations and communities* (3rd ed.). Oxford, UK: Blackwell Science Inc. <https://doi.org/10.1002/9781444313765>
- Benito, B. M., Cayuela, L., & Albuquerque, F. S. (2013). The impact of modelling choices in the predictive performance of richness maps derived from species-distribution models: Guidelines to build better diversity models. *Methods in Ecology and Evolution*, 4, 327–335. <https://doi.org/10.1111/2041-210x.12022>
- Blois, J. L., Zarnetske, P. L., Fitzpatrick, M. C., & Finnegan, S. (2013). Climate change and the past, present, and future of biotic interactions. *Science*, 341, 499–504. <https://doi.org/10.1126/science.1237184>
- Breiman, L. (2001). Random forests. *Machine Learning*, 45, 5–32. <https://doi.org/10.1023/A:1010933404324>
- Breiner, F. T., Guisan, A., Bergamini, A., & Nobis, M. P. (2015). Overcoming limitations of modelling rare species by using ensembles of small models. *Methods in Ecology and Evolution*, 6, 1210–1218. <https://doi.org/10.1111/2041-210x.12403>
- Calabrese, J. M., Certain, G., Kraan, C., & Dormann, C. F. (2014). Stacking species distribution models and adjusting bias by linking them to macroecological models. *Global Ecology and Biogeography*, 23, 99–112. <https://doi.org/10.1111/Geb.12102>
- Cohen, J. (1968). Weighted kappa: Nominal scale agreement with provision for scaled disagreement or partial credit. *The American Psychological Association*, 70, 213–220.
- Cola, V. D., Broennimann, O., Petitpierre, B., Breiner, F. T., D'Amen, M., Randin, C., ... Dubuis, A. (2016). *ecospat*: An R package to support spatial analyses and modeling of species niches and distributions. *Ecography*, 40, 774–787.
- Cord, A. F., Klein, D., Gernandt, D. S., la Rosa, J. A. P., & Dech, S. (2014). Remote sensing data can improve predictions of species richness by stacked species distribution models: A case study for Mexican pines. *Journal of Biogeography*, 41, 736–748. <https://doi.org/10.1111/jbi.12225>
- Currie, D. J., Mittelbach, G. G., Cornell, H. V., Field, R., Guégan, J. F., Hawkins, B. A., & O'Brien, E. (2004). Predictions and tests of climate-based hypotheses of broad-scale variation in taxonomic richness. *Ecology Letters*, 7, 1121–1134. <https://doi.org/10.1111/j.1461-0248.2004.00671.x>
- D'Amen, M., Bombi, P., Pearman, P. B., Schmatz, D. R., Zimmermann, N. E., & Bologna, M. A. (2011). Will climate change reduce the efficacy of protected areas for amphibian conservation in Italy? *Biological Conservation*, 144, 989–997. <https://doi.org/10.1016/j.biocon.2010.11.004>
- D'Amen, M., Dubuis, A., Fernandes, R. F., Pottier, J., Pellissier, L., & Guisan, A. (2015). Using species richness and functional traits predictions to constrain assemblage predictions from stacked species distribution models. *Journal of Biogeography*, 42, 1255–1266. <https://doi.org/10.1111/jbi.12485>
- D'Amen, M., Mateo, R. G., Pottier, J., Thuiller, W., Maiorano, L., Pellissier, L., ... Guisan, A. (2017). Improving spatial predictions of taxonomic, functional and phylogenetic diversity. *Journal of Ecology*, 106, 76–86.
- D'Amen, M., Pradervand, J. N., & Guisan, A. (2015). Predicting richness and composition in mountain insect communities at high resolution: A new test of the SESAM framework. *Global Ecology and Biogeography*, 24, 1443–1453. <https://doi.org/10.1111/geb.12357>
- D'Amen, M., Rahbek, C., Zimmermann, N. E., & Guisan, A. (2017). Spatial predictions at the community level: From current approaches to

- future frameworks. *Biological Reviews*, 92, 169–187. <https://doi.org/10.1111/brv.12222>
- Distler, T., Schuetz, J. G., Velasquez-Tibata, J., & Langham, G. M. (2015). Stacked species distribution models and macroecological models provide congruent projections of avian species richness under climate change. *Journal of Biogeography*, 42, 1–13. <https://doi.org/10.1111/jbi.12479>
- Dormann, C. F., Elith, J., Bacher, S., Buchmann, C., Carl, G., Carre, G., & Lautenbach, S. (2013). Collinearity: A review of methods to deal with it and a simulation study evaluating their performance. *Ecography*, 36, 27–46. <https://doi.org/10.1111/j.1600-0587.2012.07348.x>
- Dubuis, A., Pottier, J., Rion, V., Pellissier, L., Theurillat, J. P., & Guisan, A. (2011). Predicting spatial patterns of plant species richness: A comparison of direct macroecological and species stacking modelling approaches. *Diversity and Distributions*, 17, 1122–1131. <https://doi.org/10.1111/j.1472-4642.2011.00792.x>
- Dullinger, S., Gatringer, A., Thuiller, W., Moser, D., Zimmermann, N. E., Guisan, A., & Hulber, K. (2012). Extinction debt of high-mountain plants under twenty-first-century climate change. *Nature Climate Change*, 2, 619–622. <https://doi.org/10.1038/Nclimate1514>
- Elith, J., Graham, C. H., Anderson, R. P., Dudik, M., Ferrier, S., Guisan, A., & Zimmermann, N. E. (2006). Novel methods improve prediction of species' distributions from occurrence data. *Ecography*, 29, 129–151. <https://doi.org/10.1111/j.2006.0906-7590.04596.x>
- Elith, J., & Leathwick, J. R. (2009). Species distribution models: Ecological explanation and prediction across space and time. *Annual Review of Ecology Evolution and Systematics*, 40, 677–697. <https://doi.org/10.1146/annurev.ecolsys.110308.120159>
- Elith, J., Leathwick, J. R., & Hastie, T. (2008). A working guide to boosted regression trees. *Journal of Animal Ecology*, 77, 802–813. <https://doi.org/10.1111/j.1365-2656.2008.01390.x>
- Ferrier, S., Drielsma, M., Manion, G., & Watson, G. (2002). Extended statistical approaches to modelling spatial pattern in biodiversity in northeast New South Wales II. community-level modelling. *Biodiversity and Conservation*, 11, 2309–2338. <https://doi.org/10.1023/A:1021374009951>
- Ferrier, S., & Guisan, A. (2006). Spatial modelling of biodiversity at the community level. *Journal of Applied Ecology*, 43, 393–404. <https://doi.org/10.1111/j.1365-2664.2006.01149.x>
- Fielding, A. H., & Bell, J. F. (1997). A review of methods for the assessment of prediction errors in conservation presence-absence models. *Environmental Conservation*, 24, 38–49. <https://doi.org/10.1017/S0376892997000088>
- Fleishman, E., Noss, R. F., & Noon, B. R. (2006). Utility and limitations of species richness metrics for conservation planning. *Ecological Indicators*, 6, 543–553. <https://doi.org/10.1016/j.ecolind.2005.07.005>
- Freeman, E. A., & Moisen, G. G. (2008). A comparison of the performance of threshold criteria for binary classification in terms of predicted prevalence and kappa. *Ecological Modelling*, 217, 48–58. <https://doi.org/10.1016/j.ecolmodel.2008.05.015>
- Gotelli, N. J., Anderson, M. J., Arita, H. T., Chao, A., Colwell, R. K., Currie, D. J., & Willig, M. R. (2009). Patterns and causes of species richness: A general simulation model for macroecology. *Ecology Letters*, 12, 873–886. <https://doi.org/10.1111/j.1461-0248.2009.01353.x>
- Granger, V., Bez, N., Fromentin, J. M., Meynard, C., Jadaud, A., & Merigot, B. (2015). Mapping diversity indices: Not a trivial issue. *Methods in Ecology and Evolution*, 6, 688–696. <https://doi.org/10.1111/2041-210X.12357>
- Guisan, A., Dubuis, A., Pellissier, L., Pradervand, J. N., Meier, S., Scherrer, D., ... Vittoz, P. (2018). Data from: How to best threshold and validate stacked species assemblages? Community optimisation might hold the answer. *Dryad Digital Repository*, <https://doi.org/10.5061/dryad.nf925> ps
- Guisan, A., Dubuis, A., & Vittoz, P. (2011). Data from: Predicting spatial patterns of plant species richness: A comparison of direct macroecological and species stacking modelling approaches. *Dryad Digital Repository*, <https://doi.org/10.5061/dryad.28d4k>
- Guisan, A., Lehmann, A., Ferrier, S., Aspinall, R., Overton, R., Austin, M., & Hastie, T. (2006). Making better biogeographic predictions of species distribution. *Journal of Applied Ecology*, 43, 386–392. <https://doi.org/10.1111/j.1365-2664.2006.01164.x>
- Guisan, A., & Rahbek, C. (2011). SESAM - a new framework integrating macroecological and species distribution models for predicting spatio-temporal patterns of species assemblages. *Journal of Biogeography*, 38, 1433–1444. <https://doi.org/10.1111/j.1365-2699.2011.02550.x>
- Guisan, A., & Thuiller, W. (2005). Predicting species distribution: Offering more than simple habitat models. *Ecology Letters*, 8, 993–1009. <https://doi.org/10.1111/j.1461-0248.2005.00792.x>
- Guisan, A., Thuiller, W., & Zimmermann, N. E. (2017). *Habitat suitability and distribution models*. Cambridge, UK: Cambridge University Press. <https://doi.org/10.1017/9781139028271>
- Guisan, A., Tingley, R., Baumgartner, J. B., Naujokaitis-Lewis, I., Sutcliffe, P. R., Tulloch, A. I. T., & Buckley, Y. M. (2013). Predicting species distributions for conservation decisions. *Ecology Letters*, 16, 1424–1435. <https://doi.org/10.1111/ele.12189>
- Hartmann, P., Fouvry, P., & Horisberger, D. (2009). L'Observatoire de l'écosystème forestier du canton de Vaud: Espace de recherche appliquée | The Forest Ecosystem Observatory in Canton Vaud: A field of applied research. *Schweizerische Zeitschrift für Forstwesen*, 160, s2–s6. <https://doi.org/10.3188/szf.2009.s0002>
- Hastie, T. J., & Tibshirani, R. (1990). *Generalized additive models*. London, UK: Chapman & Hall.
- Hawkins, B. A., & Porter, E. E. (2003). Water–energy balance and the geographic pattern of species richness of western Palearctic butterflies. *Ecological Entomology*, 28, 678–686. <https://doi.org/10.1111/j.1365-2311.2003.00551.x>
- Hespanhol, H., Cezón, K., Felicísimo, Á. M., Muñoz, J., & Mateo, R. G. (2015). How to describe species richness patterns for bryophyte conservation? *Ecology and Evolution*, 5, 5443–5455. <https://doi.org/10.1002/ece3.1796>
- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. New York, NY: Springer. <https://doi.org/10.1007/978-1-4614-7138-7>
- Jimenez-Valverde, A., & Lobo, J. M. (2007). Threshold criteria for conversion of probability of species presence to either-or presence-absence. *Acta Oecologica-International Journal of Ecology*, 31, 361–369. <https://doi.org/10.1016/j.actao.2007.02.001>
- Kearney, M., & Porter, W. (2009). Mechanistic niche modelling: Combining physiological and spatial data to predict species' ranges. *Ecology Letters*, 12, 334–350. <https://doi.org/10.1111/j.1461-0248.2008.01277.x>
- Liu, C. R., Berry, P. M., Dawson, T. P., & Pearson, R. G. (2005). Selecting thresholds of occurrence in the prediction of species distributions. *Ecography*, 28, 385–393. <https://doi.org/10.1111/j.0906-7590.2005.03957.x>
- Liu, C., White, M., & Newell, G. (2013). Selecting thresholds for the prediction of species occurrence with presence-only data. *Journal of Biogeography*, 40, 778–789. <https://doi.org/10.1111/jbi.12058>
- Lomba, A., Pellissier, L., Randin, C., Vicente, J., Moreira, F., Honrado, J., & Guisan, A. (2010). Overcoming the rare species modelling paradox: A novel hierarchical framework applied to an Iberian endemic plant. *Biological Conservation*, 143, 2647–2657. <https://doi.org/10.1016/j.biocon.2010.07.007>
- Maggini, R., Lehmann, A., Zimmermann, N. E., & Guisan, A. (2006). Improving generalized regression analysis for the spatial prediction of forest communities. *Journal of Biogeography*, 33, 1729–1749. <https://doi.org/10.1111/j.1365-2699.2006.01465.x>

- Magurran, A. E., & Henderson, P. A. (2003). Explaining the excess of rare species in natural species abundance distributions. *Nature*, 422, 714–716. <https://doi.org/10.1038/nature01547>
- Manel, S., Williams, H. C., & Ormerod, S. J. (2001). Evaluating presence-absence models in ecology: The need to account for prevalence. *Journal of Applied Ecology*, 38, 921–931. <https://doi.org/10.1046/j.1365-2664.2001.00647.x>
- Marmion, M., Parviainen, M., Luoto, M., Heikkinen, R. K., & Thuiller, W. (2009). Evaluation of consensus methods in predictive species distribution modelling. *Diversity and Distributions*, 15, 59–69. <https://doi.org/10.1111/j.1472-4642.2008.00491.x>
- Mateo, R. G., de la Estrella, M., Felicísimo, A. M., Muñoz, J., & Guisan, A. (2013). A new spin on a compositionalist predictive modelling framework for conservation planning: A tropical case study in Ecuador. *Biological Conservation*, 160, 150–161. <https://doi.org/10.1016/j.biocon.2013.01.014>
- Mateo, R. G., Felicísimo, A. M., Pottier, J., Guisan, A., & Muñoz, J. (2012). Do stacked species distribution models reflect altitudinal diversity patterns? *PLoS ONE*, 7, e32586. <https://doi.org/10.1371/journal.pone.0032586>
- Mateo, R. G., Mokany, K., & Guisan, A. (2017). Biodiversity models: What if unsaturation is the rule? *Trends in Ecology & Evolution*, 32, 556–566. <https://doi.org/10.1016/j.tree.2017.05.003>
- McCullagh, P., & Nelder, J. A. (1989). *Generalized linear models* (2nd ed.). London, UK: Chapman and Hall. <https://doi.org/10.1007/978-1-4899-3242-6>
- McPherson, J. M., Jetz, W., & Rogers, D. J. (2004). The effects of species' range sizes on the accuracy of distribution models: Ecological phenomenon or statistical artefact? *Journal of Applied Ecology*, 41, 811–823. <https://doi.org/10.1111/j.0021-8901.2004.00943.x>
- Meier, E. S., Kienast, F., Pearman, P. B., Svenning, J. C., Thuiller, W., Araújo, M. B., & Zimmermann, N. E. (2010). Biotic and abiotic variables show little redundancy in explaining tree species distributions. *Ecography*, 33, 1038–1048. <https://doi.org/10.1111/j.1600-0587.2010.06229.x>
- Merow, C., Smith, M. J., Edwards, T. C., Guisan, A., McMahon, S. M., Normand, S., & Elith, J. (2014). What do we gain from simplicity versus complexity in species distribution models? *Ecography*, 37, 1267–1281. <https://doi.org/10.1111/ecog.00845>
- Mod, H. K., le Roux, P. C., Guisan, A., & Luoto, M. (2015). Biotic interactions boost spatial models of species richness. *Ecography*, 38, 913–921. <https://doi.org/10.1111/ecog.01129>
- Mokany, K., & Ferrier, S. (2011). Predicting impacts of climate change on biodiversity: A role for semi-mechanistic community-level modelling. *Diversity and Distributions*, 17, 374–380. <https://doi.org/10.1111/j.1472-4642.2010.00735.x>
- Nenzen, H. K., & Araújo, M. B. (2011). Choice of threshold alters projections of species range shifts under climate change. *Ecological Modelling*, 222, 3346–3354. <https://doi.org/10.1016/j.ecolmodel.2011.07.011>
- Nogues-Bravo, D., & Rahbek, C. (2011). Communities under climate change. *Science*, 334, 1070–1071. <https://doi.org/10.1126/science.1214833>
- Parmesan, C., & Yohe, G. (2003). A globally coherent fingerprint of climate change impacts across natural systems. *Nature*, 421, 37–42. <https://doi.org/10.1038/nature01286>
- Pellissier, L., Pradervand, J.-N., Pottier, J., Dubuis, A., Maiorano, L., & Guisan, A. (2012). Climate-based empirical models show biased predictions of butterfly communities along environmental gradients. *Ecography*, 35, 684–692. <https://doi.org/10.1111/j.1600-0587.2011.07047.x>
- Pereira, H. M., Ferrier, S., Walters, M., Geller, G. N., Jongman, R. H. G., Scholes, R. J., & Wegmann, M. (2013). Essential biodiversity variables. *Science*, 339, 277–278. <https://doi.org/10.1126/science.1229931>
- Pineda, E., & Lobo, J. M. (2009). Assessing the accuracy of species distribution models to predict amphibian species richness patterns. *Journal of Animal Ecology*, 78, 182–190. <https://doi.org/10.1111/j.1365-2656.2008.01471.x>
- Pottier, J., Dubuis, A., Pellissier, L., Maiorano, L., Rossier, L., Randin, C. F., & Guisan, A. (2013). The accuracy of plant assemblage prediction from species distribution models varies along environmental gradients. *Global Ecology and Biogeography*, 22, 52–63. <https://doi.org/10.1111/j.1466-8238.2012.00790.x>
- Pouteau, R., Bayle, E., Blanchard, E., Birnbaum, P., Cassan, J. J., Hequet, V., & Vandrot, H. (2015). Accounting for the indirect area effect in stacked species distribution models to map species richness in a montane biodiversity hotspot. *Diversity and Distributions*, 21, 1329–1338. <https://doi.org/10.1111/ddi.12374>
- Pradervand, J. N., Dubuis, A., Reymond, A., Sonnay, V., Gelin, A., & Guisan, A. (2013). Quels facteurs influencent la richesse en orthoptères des Préalpes vaudoises? *Bulletin de la Société Vaudoise des Sciences Naturelles*, 93, 155–173.
- Preston, F. W. (1948). The commonness, and rarity, of species. *Ecology*, 29, 254–283. <https://doi.org/10.2307/1930989>
- Scherrer, D., Massy, S., Meier, S., Vittoz, P., & Guisan, A. (2017). Assessing and predicting shifts in mountain forest composition across 25 years of climate change. *Diversity and Distributions*, 23, 517–528. <https://doi.org/10.1111/ddi.12548>
- Sørensen, T. (1948). A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons. *Biologiske Skrifter*, 5, 1–34.
- Thuiller, W., Georges, D., & Engler, R. (2013). biomod2: Ensemble platform for species distribution modeling. *R Package Version*, 2, r560.
- Thuiller, W., Lavorel, S., Araújo, M. B., Sykes, M. T., & Prentice, I. C. (2005). Climate change threats to plant diversity in Europe. *Proceedings of the National Academy of Sciences of the United States of America*, 102, 8245–8250. <https://doi.org/10.1073/pnas.0409902102>
- Thuiller, W., Pollock, L. J., Gueguen, M., & Münkemüller, T. (2015). From species distributions to meta-communities. *Ecology Letters*, 18, 1321–1328. <https://doi.org/10.1111/ele.12526>
- Turner, J. R., Gatehouse, C. M., & Corey, C. A. (1987). Does solar energy control organic diversity? Butterflies, moths and the British climate. *Oikos*, 48, 195–205. <https://doi.org/10.2307/3565855>
- Van der Putten, W. H., Macel, M., & Visser, M. E. (2010). Predicting species distribution and abundance responses to climate change: Why it is essential to include biotic interactions across trophic levels. *Philosophical Transactions of the Royal Society B-Biological Sciences*, 365, 2025–2034. <https://doi.org/10.1098/rstb.2010.0037>
- Zimmermann, N. E., Edwards, T. C., Graham, C. H., Pearman, P. B., & Svenning, J.-C. (2010). New trends in species distribution modelling. *Ecography*, 33, 985–989. <https://doi.org/10.1111/j.1600-0587.2010.06953.x>
- Zurell, D., Zimmermann, N. E., Sattler, T., Nobis, M. P., & Schröder, B. (2016). Effects of functional traits on the prediction accuracy of species richness models. *Diversity and Distributions*, 22, 905–917. <https://doi.org/10.1111/ddi.12450>

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**How to cite this article:** Scherrer D, D'Amen M, Fernandes RF, Mateo RG, Guisan A. How to best threshold and validate stacked species assemblages? Community optimisation might hold the answer. *Methods Ecol Evol*. 2018;9:2155–2166. <https://doi.org/10.1111/2041-210X.13041>