1 **Assessing the reliability of species distribution projections in climate change research**

2

3

4 Luca Santini [1-2*], Ana Benítez-López [2-3], Luigi Maiorano [4], Mirza Čengić [1], Mark A.J. Huijbregts [1]

5

6 [1] National Research Council, Institute of Research on Terrestrial Ecosystems (CNR-IRET), Via

7 Salaria km 29.300, 00015, Monterotondo (Rome), Italy

8

9 [2] Department of Environmental Science, Institute for Wetland and Water Research, Faculty of

10 Science, Radboud University, P.O. Box 9010, NL-6500 GL, Nijmegen, The Netherlands

11

12

13 [3] Integrative Ecology Group, Estación Biológica de Doñana (EBD-CSIC), Sevilla, Spain

14

15 [4] Department of Biology and Biotechnologies, Sapienza Università di Roma, Rome, Italy

16

17 * corresponding author: luca.santini.eco@gmail.com

18

19 **Abstract**

20 <u>Aim</u>

21 Forecasting changes in species distribution under future scenarios is one of the most prolific areas

22 of application for species distribution models (SDMs). However, no consensus yet exists on the

23 reliability of such models for drawing conclusions on species distribution response to changing

24 climate. In this study we provide an overview of common modelling practices in the field and

25 assess model predictions reliability using a virtual species approach.

26

27 <u>Location</u>

28 Global

29

30 <u>Methods</u>

31 We first provide an overview of common modelling practices in the field by reviewing the papers

32 published in the last 5 years. Then, we use a virtual species approach and three commonly applied

33 SDM algorithms (GLM, MaxEnt and Random Forest) to assess the estimated (cross-validated) and

34 actual predictive performance of models parameterized with different modelling settings and

35 violations of modelling assumptions.

36

37 <u>Results</u>

38 Our literature review shows that most papers that model species distribution under climate change

39 rely on single models (65%) and small samples (< 50 presence points, 62%), use presence-only data

40 (85%), and binarize models' output to estimate range shift, contraction or expansion (74%). Our

41 virtual species approach reveals that the estimated predictive performance tends to be over-

42 optimistic compared to the real predictive performance. Further, the binarization of predicted

43 probabilities of presence reduces models' predictive ability considerably. Sample size is one of the

44 main predictors of real accuracy, but has little influence on estimated accuracy. Finally, the

45 inclusion of irrelevant predictors and the violation of modelling assumptions increases estimated

46 accuracy but decreases real accuracy of model projections, leading to biased estimates of range

47 contraction and expansion.

48

49 <u>Main conclusions</u>

50 Our study calls for extreme caution in the application and interpretation of SDMs in the context of

51 biodiversity conservation and climate change research, especially when modelling a large number

52 of species where species-specific model settings become impracticable.

53

54

55 **Keywords**

56 Area Under the Curve (AUC), bias, climate change projections, disequilibrium, geographic extent,

57 sample size, Niche modelling, spurious relationships, True Skill Statistics (TSS).

## 1. Introduction

Understanding how climate shapes species distribution and how range shifts may be driven by future climatic change is more urgent than ever. In the last thirty years, studies aimed at developing, improving and applying species distribution models (SDMs) have proliferated (Araújo et al. 2019), and forecasting changes in species distribution under future scenarios is one of the most popular areas of application for SDMs today (Thuiller et al. 2011, Schloss et al. 2012, Newbold 2018). In SDM-based climate change forecasting studies, models are trained on current data and used to predict the probability of presence under present and future conditions. Models' predictions are often binarized to assess whether a species distribution is expected to shift, contract or expand (Newbold 2018). Although many modelling techniques require presence and absence data, many models are fitted using presence-only data, i.e., contrasting presences with random pseudo-absences, or background points, that represent available conditions (Guillera-Arroita et al. 2015). The predictive performance of these models is commonly assessed by randomly splitting the dataset into training and testing, and fitting the model on the training dataset and validating it on the testing dataset using discrimination metrics such as the True Skill Statistic (TSS) or the Area Under the Curve (AUC). While several authors have warned about the challenges and uncertainties of projecting future species distribution (Dormann 2007, Peterson et al. 2018), only few studies have tested model performance with empirical data, reporting mixed results (Araujo et al. 2005, Rapacciuolo et al. 2012, Morán-Ordóñez et al. 2017, Sofaer et al. 2018).

The literature on SDMs has grown very quickly and extensively, with papers adhering to different schools of thoughts and supporting the use of one or another technique (see Norberg et al. 2019 for an overview), suggesting different validation measures (e.g. Allouche et al. 2006, Leroy et al. 2018) or approaches (e.g. testing on spatially independent data; Bahn and McGill 2013). Additionally, a number of studies made different conclusions about the minimum number of presence points needed (Stockwell and Peterson 2002, Hernandez et al. 2006, Wisz et al. 2008, van Proosdij et al. 2016), area of sampling of background points (e.g. VanDerWal et al. 2009, Anderson and Raza 2010, Elith et al. 2010, Barve et al. 2011), or choice of environmental predictors and approaches to reduce collinearity (see Fourcade et al. 2018 for an overview). This can make it challenging and disorientating for people that approach the field of SDM for the first time. The existence of modeling software that make the application of these models easier and more accessible to people with limited modelling background (e.g. MaxEnt Phillips et al. 2006), may be counterproductive, as running an SDM in one of these software may appear simpler than it is. This is particularly worrying considering that SDMs are largely used to inform conservation science (Newbold 2018, Manish and Pandit 2019).

92    Recently, a number of experts have delineated a set of best practices, and shown that many

93    studies still apply inconsistent approaches that do not adhere to the best standards (Araújo et al.

94    2019). This generates a self-perpetuating problem, because published papers create a precedent, and

95    are used to justify modelling choices in new papers. For example, while several authors argued that

96    models' predictors should be chosen considering the biology of the species (Araújo and Guisan

97    2006, Austin and Van Niel 2011), it has become a common practice to include all bioclimatic

98    variables excluding collinear variables using automatic procedures irrespective of species-specific

99    biological considerations (e.g. Manish and Pandit 2019), increasing the risk of detecting spurious

100   relationships (Synes and Osborne 2011, Fourcade et al. 2018). Worryingly, it has been shown that

101   non-biologically relevant predictors can contribute to increase the predictive ability of the models

102   (Fourcade et al. 2018), so discrimination accuracy metrics may suggest a very good model while the

103   relationships estimated do not have a biological meaning (Journé et al. 2019, Warren et al. 2020).

104   Spurious relationships become particularly problematic when the model is projected to new areas or

105   environmental scenarios (Heikkinen et al. 2012, Bahn and McGill 2013, Yackulic et al. 2013,

106   Merow et al. 2014). Similarly, methodological papers that suggest less demanding requirements can

107   become preferred and widely cited references, reinforcing the trend. For example, van Proosdij et

108   al. (2016) concluded that 14 to 25 observations may be sufficient to run species distribution models.

109   This is now often cited to justify the use of small sample sizes (e.g. Carlson et al. 2017, Chen et al.

110   2017) despite previous recommendations suggesting a minimum of 50 points (Stockwell and

111   Peterson 2002, Hernandez et al. 2006, Wisz et al. 2008).

112   The extent to which SDMs perform adequately also depends on the degree to which

113   modelling assumptions are met. SDMs are often fitted on opportunistically collected data that

114   violate the assumption of random sampling (e.g. Guillera-Arroita et al. 2015). Furthermore, the

115   present distribution of species is rarely in equilibrium with the environment, meaning that species

116   only occupy a portion of the fundamental niche, not only because of biotic (e.g. competition or

117   predation) or dispersal constraints (e.g. physical barriers, limited dispersal abilities; Soberon and

118   Peterson 2005), but also because they may have recently contracted their distribution due to human

119   influence (e.g. Varela et al. 2009, Di Marco and Santini 2015, Faurby and Svenning 2015) or

120   stochastic events. This problem has often been discussed in the literature in relation to the

121   inferences made (Varela et al. 2009, Maiorano et al. 2013, Martínez-Freiría et al. 2016, Faurby and

122   Araújo 2018). Yet, methodological papers aimed at assessing optimal settings to run species

123   distribution models typically assume ideal conditions (e.g. van Proosdij et al. 2016).

124   Models used for future projections to inform conservation need to adhere to even higher

125   standards than those used for present predictions (Sequeira et al. 2018). In fact, while a model used

126   for predicting current distribution can still provide meaningful predictions even though the inferred

127    relationships are wrong (Fourcade et al. 2018, Warren et al. 2020), relationships need to be realistic

128    in order to make meaningful predictions to different conditions. However, although guidelines for

129    transferability have been provided (Sequeira et al. 2018), it is common to validate models on

130    present data and assume they perform equally well for future predictions. While a number of studies

131    have discussed and tested the influence of multiple sources of uncertainty on the predictive

132    accuracy of SDM predictions under present conditions (e.g. Wenger and Olden 2012, Vale et al.

133    2014, Fourcade et al. 2018, Fernandes et al. 2019), to our knowledge, no study has yet tested the

134    reliability of both present and future predictions while considering the effects of different modelling

135    settings and several violations in model assumptions simultaneously.

136         In this study we first provide an overview of common practices in the field by reviewing the

137    papers published in the last 5 years. We focused on the sample size used, choice and selection of

138    environmental predictors, types of models employed, the sampling approach of background (or

139    pseudo-absence) points, and the method used for binarization of model outputs. Then, we employed

140    a virtual species approach (Zurell et al. 2010, Meynard et al. 2019) to assess the contribution of

141    different modelling settings and violation of assumptions to the predictive accuracy and projected

142    responses to climate change of SDMs for three commonly applied model algorithms (GLM,

143    MaxEnt and RandomForest). Our approach allows validating model predictions against the virtual

144    "reality", therefore estimating true model predictive accuracy. We generated 50 virtual species

145    distributions, fitted SDMs under different conditions, and assessed the discrimination ability of

146    present and future model predictions against the real distribution. We also compared this predictive

147    ability with that estimated using a cross-validation (split-plot) approach, which is the most common

148    way of assessing model discrimination accuracy in most SDM studies. We systematically assessed

149    the combined effect of 1) the number of presence points (i.e. sample size), 2) the geographic extent

150    over which background points are drawn, 3) the number of biologically relevant (i.e. true niche

151    axes) and irrelevant predictors (i.e. spurious correlates), 4) the species prevalence (proportion of

152    study area occupied by the species), 5) the sample prevalence (proportion of presences over

153    background points), 6) the proportion of niche filling (the degree to which the species is at

154    equilibrium with the environment), 7) and the spatial bias in presence points. We then assessed

155    model predictions using two common discrimination metrics: the AUC and the TSS.

156

157

158    **2 Methods**

159

160    2.1 Literature review

161    We conducted a literature review on common practices in SDM papers that projected models to a

162    different time period (past or future). We queried Web of Science and focused on papers published

163   in the last 5 years (2015-2019) to reflect the most recent trends in the field. We randomly selected

164   50 papers per year for a total of 250 papers. From each paper we extracted the following

165   information: sample size, occurrence data type (e.g. presence-only vs. presence-absence), models

166   used, the variable selection criteria, and whether probabilistic output were binarized or not. A

167   detailed description of the literature search and data extraction is presented in Supplementary

168   material Appendix 1.

169

170   2.2 Environmental variables

171   We obtained 19 bioclimatic variables from CHELSA (http://chelsa-climate.org; Karger et al. 2017)

172   at 0.1 degree resolution (~11 km) for the present and for the future (year 2050) RCP 8.5 taking the

173   median over all the General Circulation Models (GCM). We also downloaded the human footprint

174   index for 2009 from https://wcshumanfootprint.org/ (Venter et al. 2016).

175

176   2.3 Virtual species

177   We generated 50 virtual species using the 'virtualspecies' R package (Leroy et al. 2016). For each

178   virtual species, we first determined the study area by generating a random extent between 3 and 10

179   decimal degrees (~330-1100 km) in both longitude and latitude centered around a random location

180   in the globe (Fig. 1). We then selected 6 random bioclimatic variables and sampled their values

181   within the extent using 100 random points. We used the mean and standard deviation estimated for

182   the 6 bioclimatic variables to generate the niche tolerance for the virtual species.

183        We then projected the niche within the study area for present and future conditions and

184   defined the occupied area using a threshold sampled randomly between the 0.2 to 0.8 quantiles of

185   the suitability values in the study area. The threshold is meant to represent the values above which

186   the species can survive and is assumed to be present for the validation of the distribution models

187   (see section 2.5). Note that the virtual species can potentially be present outside this study area in

188   environmentally analogous conditions, but we assume that the species is either limited by dispersal,

189   absent because of biotic interactions, or its presence outside the study area is simply unknown to the

190   modeller.

191

192   2.4 Scenario settings

193        For each virtual species, we fitted species distribution models using different cross-

194   combinations of the settings presented in Table 1.  To assess the influence of sample size, we

195   sampled random presence points (10, 25, 50, 100, 250, 500 and 1000 points) within the distribution

196   area of the species. Presences were sampled randomly and not as a function of niche suitability

197   values as there is no evidence that species abundance increases with niche suitability (Dallas and

198      Hastings 2018) and observation probability is often also a function of other factors such as

199      vegetation structure or human presence. Then, to assess the influence of the geographic extent, we

200      fitted a minimum convex polygon (MCP) around presence areas to generated a buffer expressed as

201      percentage increase of the MCP, which delimited the geographic extent within which the

202      background points were sampled (0%, 100%, 500%, 5,000%, 50,000%; the latter often resulting in

203      the entire continent). We set the number of background points depending on the number of presence

204      points and the level of sample prevalence. We used three sample prevalence: 0.01, 0.1 and 1. Not all

205      background points, however, could always be sampled depending on the selected geographic extent

206      (i.e. insufficient number of cells), leading to variable sample prevalence values.

207      In each model, we used a total number of predictor variables between 3 and 12. To assess the

208      influence of biologically relevant and irrelevant predictors of species presence, we sampled none, 3,

209      or 6 relevant bioclimatic predictors (those describing the true species niche), and none, 3 or 6

210      irrelevant bioclimatic predictors (not describing the niche) from the other 13 bioclimatic variables

211      (Table 1). Combinations yielding 0 predictor variables were not considered. We tested collinearity

212      using a stepwise VIF selection for the environmental variables in the training dataset and only

213      retained variables with VIF<3 (Zuur et al. 2010), so the final number of biologically relevant or

214      irrelevant predictors could be different from multiples of 3. As a measure of model transferability,

215      we estimated the Multivariate Environmental Similarity Surface (MESS; Elith et al. 2010) between

216      the present and future set of environmental variables used in the model fitting.

217          We also considered violations of two important assumptions underlying SDMs that are

218      common in real study cases: non-equilibrium with the environment (niche filling) and non-random

219      sampling of presence points that results in a bias along an environmental gradient. Decreasing

220      proportions of niche filling were simulated by only sampling presence points below a given quantile

221      (0.33, 0.66, 1) of the human footprint index values within the study area (Table 1). This mimics a

222      scenario where a species is potentially present (given climatic conditions) and yet absent because of

223      human impact. Note that species may be in disequilibrium with the environment for different

224      reasons (e.g. biotic interactions, dispersal limitations) but the result would be similar. For simplicity

225      we restrict our analyses to the case where species are not an equilibrium because of human impact.

226          Environmental bias was simulated by randomly sampling one of the biologically relevant

227      bioclimatic predictors used in the distribution model and sampling presences only below a given

228      quantile (0.33, 0.66, 1) of the distribution of environmental values (Table 1). This represents the

229      situation where the species has only been observed under certain conditions (i.e. sampling bias

230      correlates with environmental gradients), therefore potentially biasing the estimation of the species

231      niche. When no biologically relevant variable was included, an irrelevant predictor was selected

232      instead.

233     The full set of combinations of settings in Table 1 corresponded to 7560 models; to reduce

234     the computational effort we sampled 500 model settings from the multidimensional space using a

235     conditional Latin hypercube approach (Minasny and McBratney 2006), which ensured that the

236     subset of models is representative of the real variability occurring in the original 7560 models.

237

238     2.5 Model fitting and validation

239     We used this synthetic dataset to fit three distribution model algorithms: MaxEnt (using a 'cloglog'

240     transformation and linear and quadratic feature classes), Generalized Linear Model (GLM, with a

241     stepwise model selection based on AIC including linear and quadratic terms and weights set for

242     equal sample prevalence), and Random Forest (with stratified sampling, 500 trees, and an 'mtry'

243     parameter equal to the rounded square root of the number of predictor variables). For each, we run a

244     repeated split sample cross-validation by splitting the dataset into training (80%) and testing

245     datasets (20%) 10 times. We estimated model discrimination accuracy with the Area Under the

246     Curve (AUC) and the True Skill Statistic (TSS) (Lawson et al. 2014). Then, we fitted the model

247     using the full sample, and binarized the predictions into presence-absence by using the threshold

248     that maximized TSS. We estimated contraction and expansion areas by overlaying the binary

249     predictions for the present and the future. Finally, we validated the model predictions for the

250     present, future, and areas of contraction and expansion against the virtual reality using the same

251     discrimination metrics. This validation was performed within the area of background point

252     sampling. We matched the predicted probabilities with the true presences and absences of the

253     virtual species to estimate the true AUC, and the predicted presences and absences from the

254     binarized model with the true presences and absences of the virtual species to estimate the true TSS

255     using the threshold that maximized TSS on the testing dataset. By doing this, we were able to both

256     1) estimate model discrimination accuracy mimicking a typical ecological modeller, and 2) quantify

257     the real model discrimination accuracy by comparing the model to the virtual reality.

258
259     2.6 Evaluation of model settings

260     As a post-processing step, we used a Random Forest regression to estimate the influence of

261     different modelling settings and confounding factors on the discrimination accuracy of the three

262     distribution model algorithms. We fitted a Random Forest with 1,000 trees to each species using all

263     discrimination performance metrics (TSS and AUC, both estimated and true for the present and the

264     future) and estimated changes in distribution (% of range contraction and expansion) as dependent

265     variable (one model per dependent variable), and the values of each treatment (number of

266     presences, sample prevalence, species prevalence, environmental similarity, number of relevant

267     predictors, number of irrelevant predictors; % buffer, degree of bias in sampling points, niche filling

268    proportion) as independent variables. The 'mtry' parameter in the random forest model was set to

269    the number of predictors divided by 3 (Breiman 2001).

270         We then estimated the relative importance by permutation and partial response curves for

271    each predictor per species, and then averaged all relative importance estimates and partial response

272    curves per variable across all species. The estimated relative importance values were transformed to

273    percentages (rescaled to 100) for interpretability. Confidence intervals for both variable importance

274    and partial response curves were estimated from the standard error of the mean across species

275    models.

276

277    <u>2.7 R packages</u>

278         All analyses were computed in R v. 3.5.3 (R Core Team 2018) using the packages

279    'virtualspecies' (Leroy et al. 2016), 'raster' (Hijmans and van Etten 2014), 'PresenceAbsence'

280    (Freeman and Moisen 2015), 'dismo' (Hijmans et al. 2017), 'rgeos' (Bivand and Rundel 2013),

281    'pROC' (Robin et al. 2013), 'usdm' (Naimi 2013) and 'GISTools' (Brunsdon and Chen 2014) for

282    generating virtual species and fitting species distribution models, and 'clhs' for the conditional

283    Latin hypercube sampling (Roudier 2011). We used the R package 'randomForest' (Liaw and

284    Wiener 2002) and 'ranger' for fitting Random Forest models (Wright and Ziegler 2017), 'maxnet'

285    package to fit MaxEnt (Phillips 2017), and 'pdp' for estimating the partial response curves

286    (Greenwell, Brandon 2019). The codes used for the analyses of this paper are available as part of

287    the supplementary materials.

288

289

290    **3. Results**

291

292    3.1 Common practices in SDMs

293    Among 250 papers reviewed, 92 included correlative species distribution models projected to

294    different times, and therefore were deemed relevant for our scopes (Table S1). Based on our sample

295    and using a bootstrapping approach, we estimated that the total number of papers published

296    between 2015 and 2019 that matched this criterion is 1194-1665 (95CI), indicating that we sampled

297    approximately between 5.5 and 7.7% of the total (Appendix 1).

298    Most of the papers inspected included models fitted on relatively small sample sizes (N < 50; Fig.

299    2a), with only 18.4% including minimum samples larger than 50 and 16.1% not reporting the

300    sample size used. More than 50% of the papers included all bioclimatic variables with no biological

301    justification (Fig. 2b). Among these papers, in ~50% of the cases the authors reduced the number of

302    variables using automatized approaches based on correlations or best fit to the data. A smaller

303    number of studies selected variables a priori, some of which did not provide a justification for this

304    choice (7.8%). Most of the studies used a single model (Fig. 2c,d), with MaxEnt being the most

305    common algorithm used (78.3%), followed by linear models (GLM = 30.4%; GAM = 26.1%) and

306    machine learning models (RF = 27.2%; GBM = 20.7%) (Fig. 2c). The majority of studies did not

307    include real absences but used pseudo-absences, background data, or presence-only methods (i.e.

308    climatic envelopes; 84.8%; Fig. 2e). A large proportion of papers using pseudo-absences or

309    background points did not report the area of sampling (48.7%), while others used a variety of

310    different approaches, the most common being sampling randomly across the pre-defined study area

311    (Fig. 2f). Finally, most studies binarized the continuous probability outputs based on discrimination

312    metrics (e.g. max TSS or equal sensitivity and specificity; 73.9%), almost one quarter of the studies

313    reported the continuous output (22.8%), and a small percent (3.3%) categorized probabilities into

314    multiple arbitrary categories (e.g. 0.3, 0.6 and > 0.6; Fig. 2g).

315

316    3.2 Reliability of climate change predictions

317    The three algorithms showed a consistent pattern across the two scenarios and accuracy metrics.

318    The predictive accuracy of models' predictions estimated by cross-validation was consistently

319    above the typically accepted performance thresholds (AUC=0.7, TSS=0.5; Landis and Koch 1977,

320    Swets 1988), and higher than the true predictive accuracy for present, future predictions, and

321    contraction and expansion areas (Fig. 3). However, the accuracy of binary predictions (TSS) was

322    substantially lower than that measured for continuous predictive outputs (AUC), suggesting that the

323    binarization of relative probabilities of presence decreases models' predictive ability considerably

324    (Fig. 3). Models' predictions for the future and contraction and expansion areas showed lower

325    predictive performance (Fig. 3a), especially when binarized (Fig. 3b).

326    Under optimal modelling settings (e.g. large sample size, relevant predictors, no violation of

327    assumptions regarding niche filling and unbiased sampling), models performed relatively well

328    according to AUC (Fig. S1a), but poorly when considering binary outputs (Fig. S1c). On the

329    contrary, under poor modelling settings and conditions (small sample size, irrelevant predictors,

330    violation of the main assumptions), the estimated predictive abilities remained high, but the true

331    predictive abilities dropped considerably, especially when predictions were binarized into presence-

332    absence (Fig. S1b,d). TSS and AUC were highly correlated, but while high TSS always

333    corresponded to high AUC, the opposite was not always true (Fig. S2).

334

335    3.3 Determinants of estimated predictive ability

336    The importance and effect of different factors on the estimated predictive accuracy by cross-

337    validation was qualitatively similar when using TSS or AUC (Fig. S2-S8). The most important

338  predictors of estimated predictive accuracy were species prevalence, the environmental gradient

339  sampled (inverse of environmental bias), and the geographic extent of background point sampling

340  (Fig. 4). Additionally, sample prevalence was important for Random Forest, and the number of

341  presence points for GLM (Fig. 4). Predictive accuracy decreased with increasing species prevalence

342  and decreasing environmental gradient sampled (i.e. increased with environmental bias), and

343  increased with increasing geographic extent sampled (Fig. S3-S8). The number of presence points

344  had a positive effect when fitting GLM and MaxEnt models, but had little effect when using

345  Random Forests. Sample prevalence a had positive effect in Random Forests, and weakly negative

346  in the other two models. The number of relevant and irrelevant predictors had a weak but positive

347  effect regardless of the model (Fig. S3-S8).

348

349  3.4 Determinants of true predictive ability

350  The true predictive accuracy (i.e. measured against the virtual reality) of the models for the present

351  was mostly affected by species prevalence, the number of presences, and the environmental

352  gradient sampled. The geographic extent was also important when fitting MaxEnt models (Fig. 4).

353  Both the number of presence points and the environmental gradient sampled had a positive

354  influence on predictive accuracy, geographic extent had weak positive effect, and the species

355  prevalence a negative effect (Fig. S3-S8).

356  The number of biologically relevant and irrelevant predictors, and niche filling, were relevant for

357  present predictions, but became especially influential for the predictive accuracy of models

358  projected into the future, with the number of relevant predictors and niche filling increasing

359  predictive performance, and the number of irrelevant predictors decreasing predictive performance

360  (Fig. S3-S8). An important predictor of the predictive accuracy of future projections was the degree

361  of environmental similarity between the present and future environmental conditions of the study

362  area (Fig. 4, Fig. S3-S8).

363        Species with high prevalence were more likely to expand and less likely to contract the

364  range. However, a number of additional factors contributed to these estimates (Fig. S9-S12), such as

365  the number biologically relevant and irrelevant predictors, showing a positive effect on contraction

366  and expansion estimates in GLM and MaxEnt, and a negative effect on contraction areas in random

367  forest (Fig. S10-S12). The environmental similarity between present and future conditions yielded a

368  negative effect on contraction and expansion areas, but showed non-linearity for contraction areas

369  estimated by GLM and Random Forest models. Violation of equilibrium and random sampling

370  assumption also contributed to increase range contraction and expansion estimates (Fig. S10-S12).

371

372

**4. Discussion**

In this paper we report on common practices in SDM and use this information to assess the effects of these practices on the predictive accuracy of SDMs, and thus, on the reliability of future climate-induced range shifts. Our literature review points out that a large part of papers that model species distribution under climate change rely on single models (typically MaxEnt), include models fitted on very small samples, use presence-only data, and typically binarize models' output to measure range shift, contraction or expansion. Consistently with previous analyses (Araújo et al. 2019), it also highlighted how poor modelling practices are common in the literature, especially in relation to the use of very small samples, lack of ecological considerations in the selection of model predictors, and non-reporting of fundamental information on background sample selection and study area (Zurell et al. 2020). When exploring the influence of these practices on the predictive accuracy using a virtual species approach, we found out that the estimated discrimination capacity by TSS and AUC does not reflect the actual predictive ability of SDMs, and tends to be over-optimistic compared to the real model performance when predicted under present conditions, and especially when projected to future (different) conditions. The ability of models to discriminate presences from absences as measured by the TSS is particularly low, even under optimal model settings, good ecological knowledge of the species climatic requirements, and modelling assumptions are fully met. The extent to which predictions are reliable depends on a number of model parameters (e.g. number of presence points), actual proportion of species distribution within the geographic extent (species prevalence), our degree of knowledge of the species ecology (predictor variables included in the model), and difference between present and future environmental conditions. Under optimal settings and a good ecological knowledge of the species climatic requirements, future predictions show low discrimination ability, whereas under non-optimal settings, predictions may not be better than random. Ultimately, our results suggest that irrespective of the estimated performance, we may be unable to make meaningful future predictions for many species, and even when we can, binarization of models' outputs should be avoided. Based on our results, we elaborate in the following paragraphs on guidelines and recommendations for good modelling practices when fitting SDMs.

4.1 Aim for large sample sizes

An important determinant of predictive accuracy that is often undervalued is sample size. Previous studies suggested a minimum of 50 points (Stockwell and Peterson 2002, Hernandez et al. 2006, Wisz et al. 2008), and van Proosdij et al. (2016) suggested even fewer were needed. However, these studies assessed the number of points needed under optimal conditions where the modeller uses biologically relevant environmental predictors, points are sampled randomly, and

408  species are in equilibrium with the environment; or used real species (therefore estimating accuracy

409  with testing data, e.g. Wisz et al. 2008). Our results show that while sample size has a little

410  influence on the estimated (cross-validated) accuracy, it is one of the most important predictors of

411  true accuracy. The relationship with sample size is asymptotic, and tends to stabilize around 200-

412  500 points. We must stress, however, that no magic number exists, and these values are contingent

413  on the settings in our simulation (e.g. the number of predictor variables used in the models).

414  Because we are rarely aware if the predictor variables are directly linked to species ecology, or if

415  the species is in equilibrium with the environment or presence points are biased, one should always

416  aim for the largest possible sample. This may be impracticable for many species, that are either

417  poorly known, or narrow ranged. In the absence of biological information on e.g. species' thermal

418  tolerance, it is hard to say, however, if species that are narrow ranged are specialist of specific

419  climate conditions, or are in disequilibrium with the environment. This second case likely would

420  result in an under-estimation of niche tolerance and over-prediction of range contraction under

421  climate change (Araújo and Pearson 2005, Martínez-Freiría et al. 2016, Faurby and Araújo 2018).

422  In these cases, alternative conservation assessments should be considered when possible. Projecting

423  SDMs trained on insufficient samples does not improve our knowledge in any meaningful way and

424  may actually be detrimental.

425

426  4.2 Behold sample prevalence, not the absolute number of background points

427  Many SDM studies using presence-only data sample a large number of background points or

428  pseudo-absences (e.g. 10,000), often citing Barbet-Massin et al. (2012) or Phillips and Dudík (2008)

429  as supporting reference. However, Barbet-Massin et al. did not test MaxEnt, and showed important

430  differences between algorithms. In turn, Phillips and Dudík (2008) tested MaxEnt but they report

431  their results for many species with different numbers of presence points. Hence, the positive

432  relationship between AUC and the number of background points they found should be interpreted

433  carefully as it is mediated by sample prevalence. A recent study concluded that the number of

434  background points depends on the modelling technique used (Liu et al. 2019), with accuracy in

435  MaxEnt stabilizing above a few hundreds of background points, and large numbers being only

436  relevant for common species with small samples of training presences. Our results show that GLM

437  and MaxEnt work best when sample prevalence is very low, supporting the practice of sampling a

438  large number of background points or pseudo-absences compared to the number of presences.

439  However, matching the findings by Barbet-Massin et al. (2012), we found that Random Forest

440  models perform best with high sample prevalence. This reinforces the notion that no rule of thumb

441  exists and settings should be model- and sample-specific, which is often ignored in ensemble

442  forecasting approaches that fit all models on the same dataset (e.g. Avalos and Hernández 2015,

443  Sales et al. 2017).

444

445  <u>4.3 Choose predictors carefully</u>

446  The number and quality of predictors does not seem to have a clear effect on estimated accuracy, if

447  any, increasing the number irrespective of the true underlying relationship, tends to deceitfully

448  increase estimated performance, and increase or decrease the estimated range contraction and

449  expansion. Choosing biologically meaningful predictors may not be particularly problematic when

450  predicting to present conditions (Fourcade et al. 2018), but it becomes a serious issue when the

451  model is transferred in space or time (Wenger and Olden 2012, Sequeira et al. 2018). Here we

452  considered an optimistic scenario where only 6 climatic variables influence species distributions. In

453  reality, there might be many biologically relevant variables that determine or influence the

454  distribution of a species, but our results suggest that when model is projected under different

455  conditions is better to aim for few variables for which we have clear biological expectations than

456  many variables with unclear effects on the species' distribution (Araújo and Guisan 2006, Austin

457  and Van Niel 2011).

458

459  <u>4.4 Geographic extents should accommodate the purpose of the study</u>

460      Previous studies suggest sampling background points in areas that are potentially accessible to

461  the species (e.g. biome or continent) (Araújo et al. 2019) or considering the historical biogeography

462  of the species (Barve et al. 2011, Merow et al. 2013, Cooper and Soberón 2018). This is meant to

463  allow a fair comparison between what is used and what is available. Sampling over large areas tend

464  to inflate estimated predictive accuracy, whereas the effect on true predictive accuracy of present

465  and future predictions is inconsistent across metrics (positive for AUC and negative or flat for TSS)

466  and models. This suggests that the most appropriate geographic area for sampling background

467  points varies across species and it should be tailored to the objective of the study. Setting a

468  biologically meaningful sampling area requires a deep knowledge of species ecology (e.g. dispersal

469  distance, physical and biotic barriers) and biogeography (e.g. historical distribution), which is

470  unavailable for most species, an important future avenue of research may be delineating rules of

471  thumbs that tend to improve accuracy.

472

473  <u>4.5 Noise is inevitable</u>

474  An important driver of the variation in model performance is species prevalence (Leroy et al. 2018).

475  Our results concur with previous studies showing that generalist species are harder to predict than

476  specialist species (Evangelista et al. 2008). However, "generalist" and "specialist" are relative terms

477   in the context of species distribution models, as they are defined based on the geographic extent

478   being sampled. Species prevalence over the geographic extent is something we are unaware of in

479   real study cases, and will always be an unknown factor that affects our predictive ability (Leroy et

480   al. 2018). In this sense, we should aim to optimize other model settings that can be controlled for,

481   such as the choice of predictors, the sample prevalence or having a biologically plausible

482   geographic extent.

483   Our results also show that when future environmental conditions are very dissimilar from

484   present conditions, model's projection tend to perform poorly. While entirely expected as model

485   predictions extrapolate beyond the model domain (Elith et al. 2010), this is in a way paradoxical. In

486   fact, the more dissimilar future conditions will be, the more species are expected to shift their

487   distribution range and projections becomes important to inform conservation science. Our results

488   not only corroborate previous studies emphasizing the importance of identifying extrapolation areas

489   for highlighting projection uncertainty (Elith et al. 2010, Owens et al. 2013), but also indicate that

490   forecasting accuracy decreases substantially an already low predictive performance.

491

492   4.6 Violation of modelling assumptions provides a false sense of accuracy

493   Species distribution models rely on the assumptions of random sampling and species equilibrium

494   with the environment. Worryingly, our results show that when these two assumptions are not met,

495   the estimated accuracy by cross-validation can be inflated, therefore giving the false impression that

496   the model performs well. The extensive use of citizen science data in SDMs make models

497   particularly prone to sample bias, with points more often collected in areas highly accessible to

498   humans (Bean et al. 2012), or in countries that upload their data to platforms like GBIF more

499   consistently (Meyer et al. 2016). Bias can be controlled via a number of techniques, such as

500   including covariates that act as a proxy for the bias (Warton et al. 2013), manipulating background

501   points (Ranc et al. 2016, Vollering et al. 2019), thinning presence points across the geographic

502   (Veloz 2009) or the environmental space (de Oliveira et al. 2014), or weighting data points (Elith et

503   al. 2010). While sampling bias can be sometimes obvious when we compare our sample to the

504   known approximate distribution of the species (e.g. by using IUCN range maps, or atlases), niche

505   filling is harder to evaluate, as we only have good knowledge of the historical biogeography of a

506   relatively small number of species. In many cases, the current distribution of species may result in a

507   circular reasoning, where small ranges may suggest narrow climatic tolerance while the species

508   only persists in a given geographic area for different reasons, e.g. because of anthropogenic impact

509   (Di Marco and Santini 2016). Multiple studies have discussed the under-estimation of the niche due

510   to historical range contractions (Varela et al., 2009; Maiorano et al., 2013; Martínez-Freiría et al.,

511   2016; Faurby & Araújo, 2018), and demonstrated these may largely influence our future projections

512   (Martínez-Freiría et al. 2016, Faurby and Araújo 2018). A possibility to alleviate this effect is using

513   a multi-temporal approach (or time-calibrated models) by including historical data associated with

514   the corresponding temporal climatic variables in the model training (Nogués-Bravo 2009, Maiorano

515   et al. 2013). Yet, historical records are rarely available, so we should expect that the niche always

516   tends to be under-estimated by an unknown extent compared to the true species potential, and

517   climatic projections may therefore tend to be pessimistic on average about future species occurrence

518   (Martínez-Freiría et al. 2016, Faurby and Araújo 2018).

519

520   4.7 Binarization

521   Accuracy metrics can fool us easily, and should not be used acritically to assess the reliability of a

522   model, especially considering that they can provide higher estimates in sub-optimal conditions as

523   we have shown here (Fig. S1). Some of the problems discussed above arise from the binarization of

524   probabilistic model outputs into suitable and unsuitable areas (e.g. to determine the area of range

525   contraction or expansion) based on a threshold. In fact, the true AUC tends to perform better than

526   true TSS (Fig. S1-S2), and the estimated AUC has similar values to those of the true AUC under

527   optimal conditions, whereas the estimated TSS is consistently higher than the true TSS, thus

528   overestimating accuracy. Studies using MaxEnt typically only show AUC values (standard output

529   of the software), even though model predictions are binarized. Here we show that while AUC is, as

530   expected, highly correlated with TSS, high AUC can correspond to low TSS (Fig. S2). The problem

531   arises from the fact that even when the model performs well, the threshold that maximizes

532   discriminatory ability on the training dataset may not discriminate well true presence/absence,

533   especially under different environmental conditions. Additionally, classical cross-validation is

534   performed by using a split-sample approach, but a better and more informative option is to cross-

535   validate on spatially independent samples (Bahn and McGill 2013, Roberts et al. 2017).

536   Additionally, other performance metrics focusing on probabilities (e.g. Boyce index) should be

537   considered when possible. Several authors have argued that binarization should be entirely avoided

538   unless it is clearly justified by the model application's objective (Guillera-Arroita et al. 2015). Our

539   results support this recommendation, and actually indicate that binary outputs should never be

540   considered or used to quantify changes in distribution areas. Alternative approaches to summarize

541   the results should be considered, such as looking at trends in predicted probabilities per areas.

542

543   4.8 Additional sources of uncertainty to be considered

544   In this study we evaluated the sensitivity of SDM predictions to a number of modelling settings and

545   common violations of SDM assumptions. However, there are additional factors that we did not

546   consider that can further contribute to making predictions less reliable. These include the spatial

547  accuracy of data points in relation to the resolution used (Graham et al. 2008), the taxonomic

548  accuracy of the data points (i.e., species confused with others, especially from citizen science data),

549  and the ambiguous taxonomy of the species that may lead to merging data for different species, or

550  viceversa not considering part of the distribution of a species (Araújo et al. 2019). Furthermore,

551  species distribution models assume that the species niche is static, thereby ignoring intraspecific

552  variation and local adaptations across populations (Pearman et al. 2010, Valladares et al. 2014).

553  This can be particularly problematic in climate change studies, as populations can adapt to climate

554  change (Hoffmann and Sgró 2011), and different populations can hold diverse degrees of adaptation

555  potential (Razgour et al. 2019).

556

557  4.9 Concluding remarks

558  Estimating the distribution of a species is a non-trivial task, as it requires a careful consideration of

559  the biology of the species and its historical biogeography. Uncertainty is expected to be particularly

560  high in studies modelling hundreds or thousands of species (Warren et al. 2013, 2018, Visconti et

561  al. 2016, Newbold 2018, Thuiller et al. 2019), where species-specific considerations on the

562  geographic extent or variables to include become impracticable, and normally the same geographic

563  extent for sampling background points or set of variables is used. These studies are powerful for

564  communicating important messages at the level of geographic areas (e.g., biomes) and entire

565  communities, but need to be interpreted with extreme caution, and are ill-suited for drawing

566  inferences at the level of species.

567  Our study indicates that our ability to predict future species distribution is low under on average,

568  and can be low to the point of not being meaningful when conditions are far from optimal,

569  especially when models' predictions are binarized. Hence, SDM based climate change forecasting

570  must adhere to the highest standards, must be clearly described (Zurell et al. 2020), and the

571  estimated accuracy of models should be interpreted with extreme care, as well as the results,

572  especially in relation to the quantification of range shifts, contraction and expansion, and the

573  identification of areas that will be lost or gained. These considerations are also valid (and perhaps

574  more problematic considering the wide temporal window and static niche assumption) in the case of

575  hind-casting to paleoclimates, which is now common in studies focused on refugia and

576  phylogeography (e.g. Svenning et al. 2011). Future research may focus on developing novel

577  approaches to improve, synthesize and communicate SDM projections.

578

579 **References**

580 Allouche, O. et al. 2006. Assessing the accuracy of species distribution models: Prevalence, kappa

581     and the true skill statistic (TSS). - J. Appl. Ecol. in press.

582 Anderson, R. P. and Raza, A. 2010. The effect of the extent of the study region on GIS models of

583     species geographic distributions and estimates of niche evolution: Preliminary tests with

584     montane rodents (genus Nephelomys) in Venezuela. - J. Biogeogr. 37: 1378–1393.

585 Araujo, M. et al. 2005. Validation of species-climate impact models under climate change. - Glob.

586     Chang. Biol. 11: 1504–1513.

587 Araújo, M. B. and Pearson, R. G. 2005. Equilibrium of species' distributions with climate. -

588     Ecography (Cop.). 28: 693–695.

589 Araújo, M. B. and Guisan, A. 2006. Five (or so) challenges for species distribution modelling. - J.

590     Biogeogr. 33: 1677–1688.

591 Araújo, M. B. et al. 2019. Standards for distribution models in biodiversity assessments. - Sci. Adv.

592     5: eaat4858.

593 Austin, M. P. and Van Niel, K. P. 2011. Improving species distribution models for climate change

594     studies: Variable selection and scale. - J. Biogeogr. 38: 1–8.

595 Avalos, V. del R. and Hernández, J. 2015. Projected distribution shifts and protected area coverage

596     of range-restricted Andean birds under climate change. - Glob. Ecol. Conserv. 4: 459–469.

597 Bahn, V. and McGill, B. J. 2013. Testing the predictive performance of distribution models. - Oikos

598     122: 321–331.

599 Barbet-Massin, M. et al. 2012. Selecting pseudo-absences for species distribution models: How,

600     where and how many? - Methods Ecol. Evol. 3: 327–338.

601 Barve, N. et al. 2011. The crucial role of the accessible area in ecological niche modeling and

602     species distribution modeling. - Ecol. Modell. 222: 1810–1819.

603 Bean, W. T. et al. 2012. The effects of small sample size and sample bias on threshold selection and

604     accuracy assessment of species distribution models. - Ecography (Cop.). 35: 250–258.

605 Bivand, R. and Rundel, C. 2013. rgeos: Interface to Geometry Engine - Open Source (GEOS). - R

606     Packag. version 0.3-2: 61.

607 Breiman, L. 2001. Machine Learning. - Mach. Learn. 45: 5–32.

608 Brunsdon, C. and Chen, H. 2014. GISTools: some further GIS capabilities for R. - R Packag.

609     version 0.7-4. https//CRAN.R-project.org/package=GISTools. in press.

610 Carlson, C. J. et al. 2017. Parasite biodiversity faces extinction and redistribution in a changing

611     climate. - Sci. Adv. 3: e1602422.

612 Chen, Y. et al. 2017. Assessing the effectiveness of China's protected areas to conserve current and

613     future amphibian diversity. - Divers. Distrib. 23: 146–157.

614   Cooper, J. C. and Soberón, J. 2018. Creating individual accessible area hypotheses improves
615       stacked species distribution model performance. - Glob. Ecol. Biogeogr. 27: 156–165.
616   Dallas, T. and Hastings, A. 2018. Habitat suitability estimated by niche models is largely unrelated
617       to species abundance. - Glob. Ecol. Biogeogr. 27: 1448–1456.
618   de Oliveira, G. et al. 2014. Evaluating, partitioning, and mapping the spatial autocorrelation
619       component in ecological niche modeling: A new approach based on environmentally
620       equidistant records. - Ecography (Cop.). 37: 637–647.
621   Di Marco, M. and Santini, L. 2015. Human pressures predict species' geographic range size better
622       than biological traits. - Glob. Chang. Biol. 21: 2169–2178.
623   Di Marco, M. and Santini, L. 2016. Climatic tolerance or geographic breadth: what are we
624       measuring? - Glob. Chang. Biol. 22: 972–973.
625   Dormann, C. F. 2007. Promising the future? Global change projections of species distributions. -
626       Basic Appl. Ecol. 8: 387–397.
627   Elith, J. et al. 2010. The art of modelling range-shifting species. - Methods Ecol. Evol. 1: 330–342.
628   Evangelista, P. H. et al. 2008. Modelling invasion for a habitat generalist and a specialist plant
629       species. - Divers. Distrib. 14: 808–817.
630   Faurby, S. and Svenning, J. C. 2015. Historic and prehistoric human-driven extinctions have
631       reshaped global mammal diversity patterns. - Divers. Distrib. 21: 1155–1166.
632   Faurby, S. and Araújo, M. B. 2018. Anthropogenic range contractions bias species climate change
633       forecasts. - Nat. Clim. Chang. 8: 252.
634   Fernandes, R. F. et al. 2019. Effects of simulated observation errors on the performance of species
635       distribution models. - Divers. Distrib. 25: 400–413.
636   Fourcade, Y. et al. 2018. Paintings predict the distribution of species, or the challenge of selecting
637       environmental predictors and evaluation statistics. - Glob. Ecol. Biogeogr. 27: 245–256.
638   Freeman, E. A. and Moisen, G. 2015. PresenceAbsence□: An R Package for Presence Absence
639       Analysis. - J. Stat. Softw. 23: 1–31.
640   Graham, C. H. et al. 2008. The influence of spatial errors in species occurrence data used in
641       distribution models. - J. Appl. Ecol. 45: 239–247.
642   Greenwell, Brandon, M. 2019. pdp: An R Package for Constructing Partial Dependence Plots. - R J.
643       9: 421–436.
644   Guillera-Arroita, G. et al. 2015. Is my species distribution model fit for purpose? Matching data and
645       models to applications. - Glob. Ecol. Biogeogr. 24: 276–292.
646   Heikkinen, R. K. et al. 2012. Does the interpolation accuracy of species distribution models come at
647       the expense of transferability? - Ecography (Cop.). 35: 276–288.
648   Hernandez, P. A. et al. 2006. The effect of sample size and species characteristics on performance

649      of different species distribution modeling methods. - Ecography (Cop.). 29: 773–785.

650    Hijmans, R. J. and van Etten, J. 2014. raster: Geographic data analysis and modeling. - R Packag.

651      version 2.3-12. http//CRAN.R-project.org/package=raster in press.

652    Hijmans, R. J. et al. 2017. Package ' dismo .' - R Packag. version 1.1-4. https//CRAN.R-

653      project.org/package=dismo 9: 1–68.

654    Hoffmann, A. A. and Sgró, C. M. 2011. Climate change and evolutionary adaptation. - Nature 470:

655      479–485.

656    Journé, V. et al. 2019. Correlative climatic niche models predict real and virtual species

657      distributions equally well. - Ecology in press.

658    Karger, D. N. et al. 2017. Climatologies at high resolution for the earth's land surface areas. - Sci.

659      Data 4: 170122.

660    Landis, J. R. and Koch, G. G. 1977. The Measurement of Observer Agreement for Categorical

661      Data. - Biometrics 33: 159–174.

662    Lawson, C. R. et al. 2014. Prevalence, thresholds and the performance of presence-absence models.

663      - Methods Ecol. Evol. 5: 54–64.

664    Leroy, B. et al. 2016. virtualspecies, an R package to generate virtual species distributions. -

665      Ecography (Cop.). 39: 599–607.

666    Leroy, B. et al. 2018. Without quality presence–absence data, discrimination metrics such as TSS

667      can be misleading measures of model performance. - J. Biogeogr. 45: 1994–2002.

668    Liaw, A. and Wiener, M. 2002. Classification and Regression by randomForest. - R News 2: 18–22.

669    Liu, C. et al. 2019. The effect of sample size on the accuracy of species distribution models:

670      considering both presences and pseudo-absences or background sites. - Ecography (Cop.). 42:

671      535–548.

672    Maiorano, L. et al. 2013. Building the niche through time: Using 13,000 years of data to predict the

673      effects of climate change on three tree species in Europe. - Glob. Ecol. Biogeogr. 22: 302–317.

674    Manish, K. and Pandit, M. K. 2019. Identifying conservation priorities for plant species in the

675      Himalaya in current and future climates: A case study from Sikkim Himalaya, India. - Biol.

676      Conserv. 233: 176–184.

677    Martínez-Freiría, F. et al. 2016. Contemporary niche contraction affects climate change predictions

678      for elephants and giraffes. - Divers. Distrib. 22: 432–444.

679    Merow, C. et al. 2013. A practical guide to MaxEnt for modeling species' distributions: What it

680      does, and why inputs and settings matter. - Ecography (Cop.). 36: 1058–1069.

681    Merow, C. et al. 2014. What do we gain from simplicity versus complexity in species distribution

682      models? - Ecography (Cop.). 37: 1267–1281.

683    Meyer, C. et al. 2016. Multidimensional biases, gaps and uncertainties in global plant occurrence

684      information. - Ecol. Lett. 19: 992–1006.

685      Meynard, C. N. et al. 2019. Testing methods in species distribution modelling using virtual species:

686      what have we learnt and what are we missing? - Ecography (Cop.). in press.

687      Minasny, B. and McBratney, A. B. 2006. A conditioned Latin hypercube method for sampling in

688      the presence of ancillary information. - Comput. Geosci. 32: 1378–1388.

689      Morán-Ordóñez, A. et al. 2017. Evaluating 318 continental-scale species distribution models over a

690      60-year prediction horizon: what factors influence the reliability of predictions? - Glob. Ecol.

691      Biogeogr. 26: 371–384.

692      Naimi, B. 2013. usdm: Uncertainty analysis for species distribution models. - R Packag. version

693      1.1-18. https//cran.r-project.org/web/packages/usdm/index.html: 1–15.

694      Newbold, T. 2018. Future effects of climate and land-use change on terrestrial vertebrate

695      community diversity under different scenarios. - Proc. R. Soc. B Biol. Sci. 285: 20180792.

696      Nogués-Bravo, D. 2009. Predicting the past distribution of species climatic niches. - Glob. Ecol.

697      Biogeogr. 18: 521–531.

698      Norberg, A. et al. 2019. A comprehensive evaluation of predictive performance of 33 species

699      distribution models at species and community levels. - Ecol. Monogr. in press.

700      Owens, H. L. et al. 2013. Constraints on interpretation of ecological niche models by limited

701      environmental ranges on calibration areas. - Ecol. Modell. 263: 10–18.

702      Pearman, P. B. et al. 2010. Within-taxon niche structure: Niche conservatism, divergence and

703      predicted effects of climate change. - Ecography (Cop.). 33: 990–1003.

704      Peterson, A. T. et al. 2018. Major challenges for correlational ecological niche model projections to

705      future climate conditions. - Ann. N. Y. Acad. Sci. 1429: 66–77.

706      Phillips, S. 2017. maxnet: fitting 'Maxent'species distribution models with 'glmnet.' in press.

707      Phillips, S. J. and Dudík, M. 2008. Modeling of species distributions with Maxent: New extensions

708      and a comprehensive evaluation. - Ecography (Cop.). 31: 161–175.

709      Phillips, S. J. et al. 2006. Maximum entropy modeling of species geographic distributions. - Ecol.

710      Modell. 190: 231–259.

711      R Core Team 2018. R: A language and environment for statistical computing. - R Found. Stat.

712      Comput. Vienna, Austria. https//www.R-project.org/ in press.

713      Ranc, N. et al. 2016. Performance tradeoffs in target-group bias correction for species distribution

714      models. - Ecography (Cop.). 40: 1076–1087.

715      Rapacciuolo, G. et al. 2012. Climatic associations of British species distributions show good

716      transferability in time but low predictive accuracy for range change. - PLoS One 7: e40212.

717      Razgour, O. et al. 2019. Considering adaptive genetic variation in climate change vulnerability

718      assessment reduces species range loss projections. - Proc. Natl. Acad. Sci. U. S. A. 116:

719      10418–10423.

720    Roberts, D. R. et al. 2017. Cross-validation strategies for data with temporal, spatial, hierarchical, or

721      phylogenetic structure. - Ecography (Cop.). 40: 913–929.

722    Robin, X. et al. 2013. pROC: an open-source package for R and S+ to analyze and compare ROC

723      curves. - BMC Bioinformatics 12: 77.

724    Roudier, P. 2011. clhs: a R package for conditioned Latin hypercube sampling. in press.

725    Sales, L. P. et al. 2017. Model uncertainties do not affect observed patterns of species richness in

726      the Amazon. - PLoS One in press.

727    Schloss, C. A. et al. 2012. Dispersal will limit ability of mammals to track climate change in the

728      Western Hemisphere. - Proc. Natl. Acad. Sci. USA 109: 8606–8611.

729    Sequeira, A. M. M. et al. 2018. Transferring biodiversity models for conservation: Opportunities

730      and challenges. - Methods Ecol. Evol. 9: 1250–1264.

731    Soberon, J. and Peterson, A. T. 2005. Interpretation of Models of Fundamental Ecological Niches

732      and Species' Distributional Areas. - Biodivers. Informatics 2: 1–10.

733    Sofaer, H. R. et al. 2018. Misleading prioritizations from modelling range shifts under climate

734      change. - Glob. Ecol. Biogeogr. 27: 658–666.

735    Stockwell, D. R. B. and Peterson, A. T. 2002. Effects of sample size on accuracy of species

736      distribution models. - Ecol. Modell. 148: 1–13.

737    Svenning, J. C. et al. 2011. Applications of species distribution modeling to paleobiology. - Quat.

738      Sci. Rev. 30: 2930–2947.

739    Swets, J. A. 1988. Measuring the accuracy of diagnostic systems. - Science (80-. ). 240: 1285–1293.

740    Synes, N. W. and Osborne, P. E. 2011. Choice of predictor variables as a source of uncertainty in

741      continental-scale species distribution modelling under climate change. - Glob. Ecol. Biogeogr.

742      20: 904–914.

743    Thuiller, W. et al. 2011. Consequences of climate change on the tree of life in Europe. - Nature 470:

744      531–534.

745    Thuiller, W. et al. 2019. Uncertainty in ensembles of global biodiversity scenarios. - Nat. Commun.

746      10: 1446.

747    Vale, C. G. et al. 2014. Predicting species distribution at range margins: Testing the effects of study

748      area extent, resolution and threshold selection in the Sahara-Sahel transition zone. - Divers.

749      Distrib. 20: 20–33.

750    Valladares, F. et al. 2014. The effects of phenotypic plasticity and local adaptation on forecasts of

751      species range shifts under climate change. - Ecol. Lett. 17: 1351–1364.

752    van Proosdij, A. S. J. et al. 2016. Minimum required number of specimen records to develop

753      accurate species distribution models. - Ecography (Cop.). 39: 542–552.

754   VanDerWal, J. et al. 2009. Abundance and the Environmental Niche: Environmental Suitability
755        Estimated from Niche Models Predicts the Upper Limit of Local Abundance. - Am. Nat. 174:
756        282–291.

757   Varela, S. et al. 2009. Is current climatic equilibrium a guarantee for the transferability of
758        distribution model predictions? A case study of the spotted hyena. - J. Biogeogr. 36: 1645–
759        1655.

760   Veloz, S. D. 2009. Spatially autocorrelated sampling falsely inflates measures of accuracy for
761        presence-only niche models. - J. Biogeogr. 36: 2290–2299.

762   Venter, O. et al. 2016. Sixteen years of change in the global terrestrial human footprint and
763        implications for biodiversity conservation. - Nat. Commun. 7: 11.

764   Visconti, P. et al. 2016. Projecting global biodiversity indicators under future development
765        scenarios. - Conserv. Lett. 9: 5–13.

766   Vollering, J. et al. 2019. Bunching up the background betters bias in species distribution models. -
767        Ecography (Cop.). in press.

768   Warren, R. et al. 2013. Quantifying the benefit of early climate change mitigation in avoiding
769        biodiversity loss. - Nat. Clim. Chang. 3: 678.

770   Warren, R. et al. 2018. The projected effect on insects, vertebrates, and plants of limiting global
771        warming to 1.5°C rather than 2°C. - Science (80-. ). 360: 791–795.

772   Warren, D. L. et al. 2020. Evaluating presence☐only species distribution models with
773        discrimination accuracy is uninformative for many applications. - J. Biogeogr. 47: 167–180.

774   Warton, D. I. et al. 2013. Model-based control of observer bias for the analysis of presence-only
775        data in ecology. - PLoS One 8: e79168.

776   Wenger, S. J. and Olden, J. D. 2012. Assessing transferability of ecological models: an
777        underappreciated aspect of statistical validation. - Methods Ecol. Evol. 3: 260–267.

778   Wisz, M. S. et al. 2008. Effects of sample size on the performance of species distribution models. -
779        Divers. Distrib. 14: 763–773.

780   Wright, M. N. and Ziegler, A. 2017. ranger☐: A Fast Implementation of Random Forests for High
781        Dimensional Data in C++ and R. - J. Stat. Softw. 77: 1–17.

782   Yackulic, C. B. et al. 2013. Presence☐only modelling using MAXENT: when can we trust the
783        inferences? - Methods Ecol. Evol. 4: 236–243.

784   Zurell, D. et al. 2010. The virtual ecologist approach: Simulating data and observers. - Oikos 119:
785        622–635.

786   Zurell, D. et al. 2020. A standard protocol for reporting species distribution models. - Ecography
787        (Cop.). in press.

788   Zuur, A. F. et al. 2010. A protocol for data exploration to avoid common statistical problems. -

789        Methods Ecol. Evol. 1: 3–14.
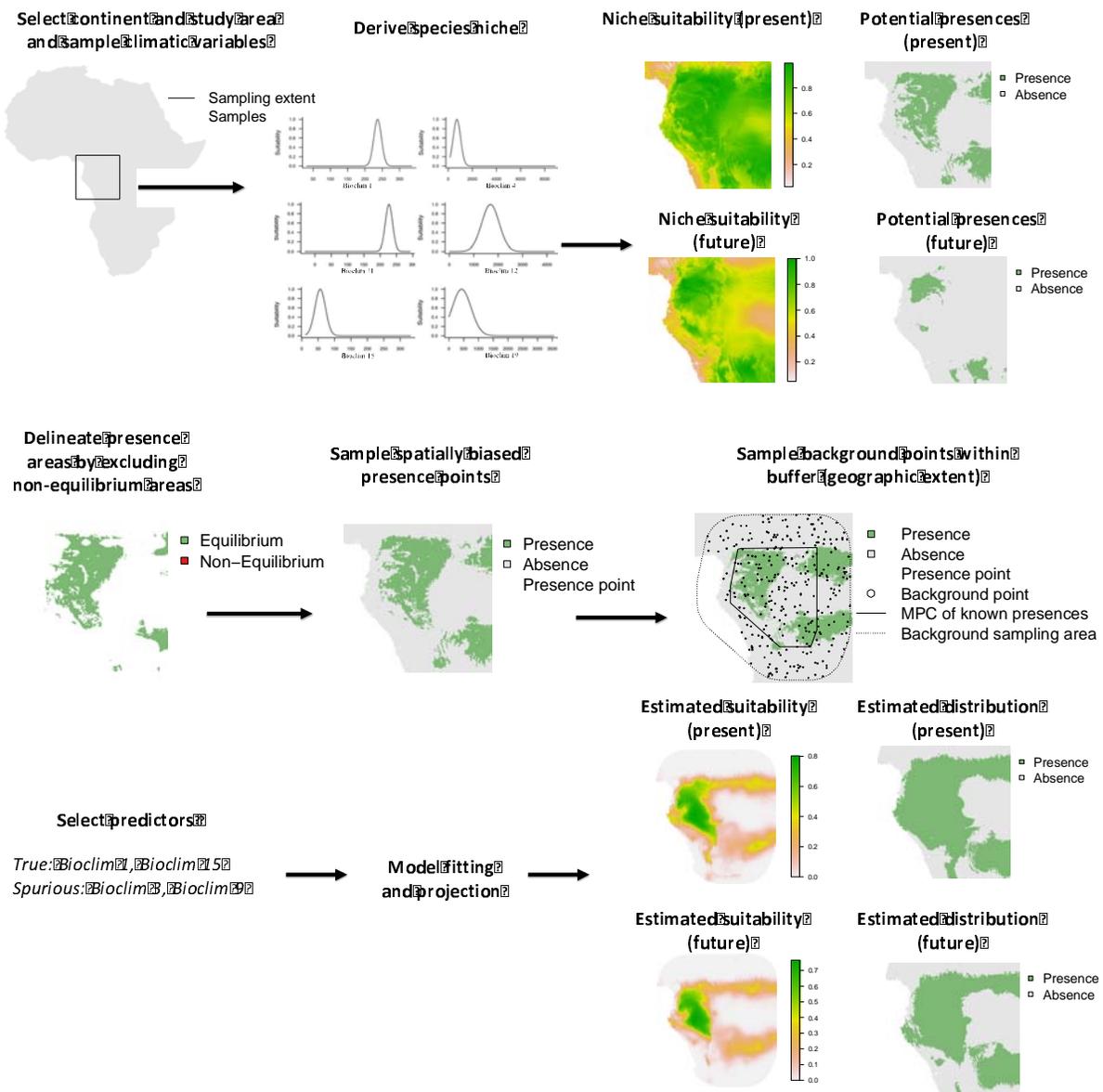
790

**Table 1.** Summary of treatments considered for fitting the species distribution models on virtual species.

| Treatment | Values |
|---|---|
| Presences | 10, 25, 50, 100, 250, 500, 1000 |
| Sample Prevalence | 0.01, 0.1, 1 |
| Buffer (%) | 0, 100, 500, 5000, 50000 |
| Bias (%) | 33, 66, 100 |
| Niche filling (%) | 33, 66, 100 |
| Relevant predictors | 0, 3, 6 |
| Irrelevant predictors | 0, 3, 6 |

794  **Fig. 1.** Modelling steps taken to generate virtual species and fit and project the species distribution
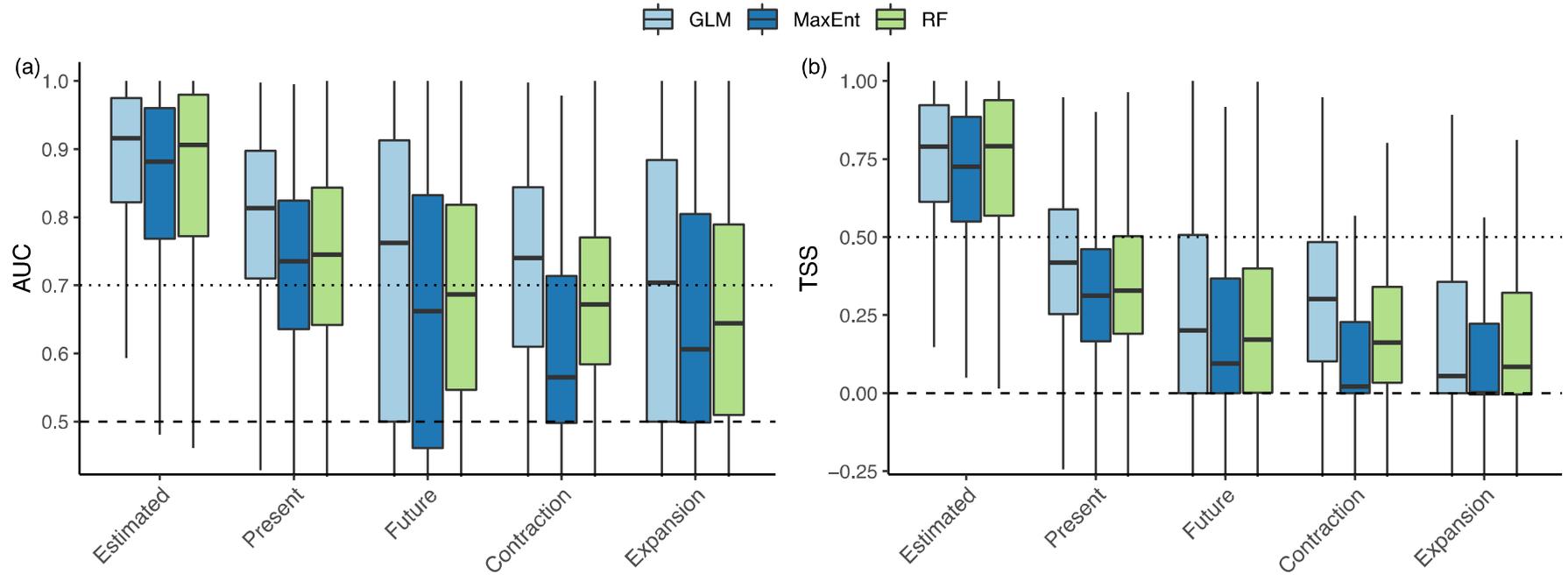795  model.
796

797 **Fig. 2.** Summary of the literature review. (a) Number of presences used in the models (minimum
798 among species if multiple species were modelled); PtP = Polygons converted to presence points;
799 NR = sample size not reported; (b) Variable selection approach. all ClmVr = all climatic variables
800 were considered; sel. ClmVr = a subset of climatic variable was considered; Automatic selection =
801 collinear variables were excluded using automatized approaches based on correlations, variance
802 inflation factors, or best fit to the data; Informed selection = collinear variables were excluded
803 based on expert opinion; No selection = Collinear variables were not excluded, or no collinearity
804 was not found or reported; No justification = no justification provided for the rationale underlying
805 the subset of variables chosen; (c) Percentages of studies using MaxEnt, generalized linear models
806 (GLM), generalized additive models (GAM), random forests (RF), generalized boosted trees
807 (GBM), multivariate adaptive regression splines (MARS), classification trees (CTA), artificial
808 neural networks (ANN) and flexible discriminant analysis (FDA). Other models used in a minority
809 of instances are not reported here (see Table S1); (d) Percentage of studies using one or multiple
810 models, or ensemble modelling approach; (e) Pseudo-absences or background points sampling
811 approach; Bias = sampling that mimics sampling bias; Buffer = random sampling within a buffer
812 around presence points; Distance-weighted = Sampling with higher intensity near (-) or far from (+)
813 presence points; Globally = Random sampling globally; Outside climate envelope = Beyond
814 climatic conditions observed for presence points; Study area = random sampling within a pre-
815 defined study area; (f) Percentage of studies using presences only, presences + background points /
816 pseudo-absences, and those using presences and real absences; (g) Percentage of studies binarizing
817 the probabilities into suitable/unsuitable, or in multiple arbitrary categories, or not applying any
818 form of binarization.

819



820

821 **Fig. 3.** AUC (a) and TSS (b) of the models fitted. Estimated = Estimated through internal cross-validation; Present and Future = True value validated
822 against virtual reality for present and future; Contraction and Expansion = True value validated against virtual reality for predicted contraction and
823 expansion areas; Dashed line = null expectation (no better than random); Dotted line = Value typically considered as "good" performance thresholds.
824 The box edges are the 25th and 75th percentiles of the distribution, and whiskers 1.5 the inter-quartile range.

825
826
827
828



829

**Fig. 4.** Relative variable importance of different settings and conditions on the TSS (a, b, c) and AUC (d, e, f) estimated by cross-validation (a, d), and measured against virtual reality for the present (b, e) and future predictions (c, f). Relative importance values are rescaled to 100 for each species. Bars represent the mean over all virtual species and error bars the standard error around the mean.