



**HAL**  
open science

# USE it: Uniformly sampling pseudo-absences within the environmental space for applications in habitat suitability models

Daniele da Re, Enrico Tordoni, Jonathan Roger Michel Henri Lenoir, Jonas Lembrechts, Sophie Vanwambeke, Duccio Rocchini, Manuele Bazzichetto

## ► To cite this version:

Daniele da Re, Enrico Tordoni, Jonathan Roger Michel Henri Lenoir, Jonas Lembrechts, Sophie Vanwambeke, et al. USE it: Uniformly sampling pseudo-absences within the environmental space for applications in habitat suitability models. *Methods in Ecology and Evolution*, 2024, 14 (11), pp.2873-2887. 10.1111/2041-210X.14209 . hal-04261748

**HAL Id: hal-04261748**

**<https://hal.science/hal-04261748v1>**

Submitted on 31 Oct 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

USE: a novel approach to uniformly sample the environmental space

1           **USE it: uniformly sampling pseudo-absences**  
2           **within the environmental space for applications**  
3           **in habitat suitability models**

4           Daniele Da Re<sup>1,\*†</sup>, Enrico Tordoni<sup>2†</sup>, Jonathan Lenoir<sup>3</sup>,

5           Jonas J. Lembrechts<sup>4</sup>, Sophie O. Vanwambeke<sup>1</sup>,

6           Duccio Rocchini<sup>5,6</sup>, and Manuele Bazzichetto<sup>7†</sup>

7           <sup>1</sup> Georges Lemaître Center for Earth and Climate Research, Earth and Life Institute,  
8           UCLouvain, Place Louis Pasteur 3, 1348 Louvain-la-Neuve, Belgium.

9           <sup>2</sup> Department of Botany, Institute of Ecology and Earth Sciences, University of Tartu, J.  
10           Liivi 2, 50409 Tartu, Estonia

11           <sup>3</sup> UMR CNRS 7058 «Ecologie et Dynamique des Systèmes Anthropisés» (EDYSAN),  
12           Université de Picardie Jules Verne, 1 rue des Louvels, 80000 Amiens, France

13           <sup>4</sup> Research Group Plants and Ecosystems, University of Antwerp, Belgium

14           <sup>5</sup> BIOME Lab., Department of Biological, Geological and Environmental Sciences, Alma  
15           Mater Studiorum University of Bologna, Via Irnerio 42, 40126 Bologna, Italy

16           <sup>6</sup> Department of Spatial Sciences, Faculty of Environmental Sciences, Czech University  
17           of Life Sciences Prague, Kamýcka 129, 16500 Praha, Czech Republic

18           <sup>7</sup> Faculty of Environmental Sciences, Department of Spatial Sciences, Czech University  
19           of Life Sciences Prague, Kamýcka 129, 16500, Praha-Suchbát, Czech Republic

20  
21           †DDR, ET and MB equally contributed to this work

22  
23           **Corresponding author:** Daniele Da Re, daniele.dare@uclouvain.be  
24  
25  
26

## 28 Abstract

- 29 1. Habitat suitability models infer the geographical distribution of species using  
30 occurrence data and environmental variables. While data on species presence are  
31 increasingly accessible, the difficulty to confirm real absences in the field often forces  
32 researchers to generate them *in silico*. To this aim, pseudo-absences are commonly  
33 randomly sampled across the study area (i.e., the geographical space). However, this  
34 introduces sample location bias (i.e., the sampling is unbalanced towards the most  
35 frequent habitats occurring within the geographical space) and favours class overlap  
36 (i.e., overlap between environmental conditions associated with species presences  
37 and pseudo-absences) in the training dataset.
- 38 2. To mitigate this, we propose an alternative methodology (i.e., the uniform approach)  
39 that systematically samples pseudo-absences within a portion of the environmental  
40 space delimited by a kernel-based filter, which seeks to minimise the number of false-  
41 absences included in the training dataset.
- 42 3. We simulated 50 virtual species and modelled their distribution using training datasets  
43 assembled with the presence points of the virtual species and pseudo-absences  
44 collected using the uniform approach and other approaches that randomly sample  
45 pseudo-absences within the geographical space. We compared the predictive  
46 performance of habitat suitability models and evaluated the extent of sample location  
47 bias and class overlap associated with the different sampling strategies.
- 48 4. Results indicated that the uniform approach: (i) effectively reduces sample location  
49 bias and class overlap; (ii) provides comparable predictive performance to sampling  
50 strategies carried out in the geographical space; and (iii) ensures gathering pseudo-

USE: a novel approach to uniformly sample the environmental space

51                    absences adequately representing the environmental conditions available across the  
52                    study area. We developed a set of R functions in an accompanying R package called  
53                    USE to disseminate the uniform approach.

54                    **Keywords:** background points, ecological niche models, presence-only models,  
55                    sample location bias, class overlap, species distribution models, reproducibility.

USE: a novel approach to uniformly sample the environmental space

## 56 1 Introduction

57 Habitat suitability models (hereafter, HSMs) are a class of statistical models used to  
58 describe the relationship between species attributes (e.g., presence-absence, abundance)  
59 and a set of spatially-explicit variables chiefly representing abiotic, biotic and human-related  
60 factors (e.g., climate, soil, demographic parameters, land-use). These models are rooted in  
61 the niche theory (i.e., *Hutchinsonian* niche, see Guisan et al., 2017) and rely on both  
62 theoretical and practical assumptions: (i) species are assumed to be at (quasi)equilibrium  
63 with their environment (Hattab et al., 2017); (ii) the set of predictors used to fit HSMs  
64 includes all necessary information to capture the ecological niche of the species; and (iii)  
65 species distribution attributes, used as the response variable, need to be appropriate for the  
66 intended model purpose (e.g., biodiversity conservation, forecasting biological invasions,  
67 assessing the effects of global change) (Tessarolo et al., 2021; but see also Guisan et al.,  
68 2017 for a thorough review on the theoretical assumptions underpinning HSMs). Some of  
69 these assumptions are hardly, if ever, met in nature since species are seldom at equilibrium  
70 with their environment (Svenning and Skov, 2004), posing several limitations to the use and  
71 interpretation of HSMs' outputs. Acknowledging and, when possible, addressing these  
72 limitations still makes HSMs a powerful toolbox for understanding the drivers of the species'  
73 realised and potential distributions (*sensu* Jackson and Overpeck, 2000). For this reason,  
74 HSMs are still widely applied in several research fields, including biogeography (Wasof et  
75 al., 2015; Duffy et al., 2017), climate change ecology (Jarvie and Svenning, 2018),  
76 conservation biology (Newbold, 2018; Santini et al., 2021), invasion ecology (Hattab et al.,  
77 2017; Da Re et al. 2020; Bazzichetto et al. 2021), and pathogen risk assessment (Batista  
78 et. al., 2023).

79 One of the most critical assumptions underpinning HSMs is the appropriateness of

USE: a novel approach to uniformly sample the environmental space

80 biological data for modelling the ecological niche of the species, which means that species  
81 distribution attributes, being either presence-absence or abundance data, should allow  
82 effectively describing the true species-environment relationship (Guisan et al., 2017; Baker  
83 et al., 2022). However, while information on species occurrence (i.e., presence) is usually  
84 readily accessible through field-collected observations or museum/herbaria records,  
85 trustworthy absence data are by far more difficult to gather or to confirm in the field  
86 (Jiménez-Valverde et al., 2008), as their sampling requires labour-intensive and costly field  
87 campaigns (Hattab et al., 2017). The usual lack of true absence data has led to the  
88 development of HSMs approaches that either rely solely on presence data (so-called  
89 'presence-only models', such as the BIOCLIM model; Booth et al. 2014) or combine  
90 presence data with pseudo-absences or background points for modelling species  
91 distributions (e.g., the MaxEnt algorithm; Phillips et al., 2017).

92 Pseudo-absences and background points are terms often used interchangeably in the  
93 scientific literature (Sillero and Barbosa, 2020), but they may represent different conditions.  
94 Pseudo-absences are sampled from locations considered unsuitable for the species  
95 (Barbet-Massin et al., 2012). In contrast, background points encompass the full range of  
96 environmental conditions, including potential suitable locations for the species (presence  
97 locations; Phillips et al., 2009; Hallgren et al., 2019). The choice between pseudo-absences  
98 and background points indicates the user's uncertainty about the ecological preferences of  
99 the species, with background points used when there is no prior knowledge of unsuitable  
100 environmental conditions. Despite recognizing the distinction, we will henceforth use the  
101 term pseudo-absences to refer to both pseudo-absences and background points for  
102 simplicity and alignment with our study.

103 The most common approaches for sampling pseudo-absences involve (i) randomly  
104 surveying a large number of points across the study area (e.g., 10,000; Barbet-Massin et

USE: a novel approach to uniformly sample the environmental space

105 al., 2012; Iturbide et al., 2015; Støa et al., 2019, Hysen et al., 2022) or (ii) sampling them  
106 within or (iii) outside buffers created around presence locations (VanDerWal et al., 2009;  
107 Bedia et al., 2013). These approaches share the characteristic of deploying pseudo-  
108 absences randomly across the geographic space, which often leads to oversampling of the  
109 most common habitat conditions that are widespread in the study area (Tessarolo et al.,  
110 2014, 2021; Ronquillo et al., 2020). This sample location bias negatively impacts HSMs in  
111 multiple ways. Firstly, it can introduce a bias in the sampling of environmental conditions  
112 experienced by a species, potentially affecting the accurate estimation of the species  
113 response curve, particularly in heterogeneous areas (Austin 2007; Hortal et al., 2008; Albert  
114 et al., 2010; Beck et al., 2014, Bazzichetto et al., 2023). Secondly, it influences the  
115 predictive performance of HSMs, as reflected in the evaluation metrics used (Jiménez-  
116 Valverde et al., 2013; Sillero and Barbosa, 2020).

117 To overcome this issue, previous studies (Varela et al. 2014; Hattab et al., 2017)  
118 proposed to sample species presence and (true) absence data throughout a systematic  
119 sampling of the environmental conditions available across the study area, thus limiting the  
120 artificial constraint imposed by the random sampling towards the most widespread  
121 environments. More specifically, Varela et al. (2014), Hattab et al. (2017) and Perret and  
122 Sax (2022) suggested collecting species' presence and/or absence within 2- or 3-  
123 dimensional environmental spaces obtained using ordination techniques. Such approaches  
124 significantly contributed to the improvement and standardisation of the way species  
125 observations, including pseudo-absences, can be collected to calibrate HSMs reducing  
126 sample location bias. Yet, they do not explicitly consider class overlap, another relevant  
127 methodological issue encountered when collecting pseudo-absences through random  
128 sampling across the geographical space. Class overlap refers to the overlap between  
129 environmental conditions associated with both species presence and absence, thus

USE: a novel approach to uniformly sample the environmental space

130 hindering the concept of pseudo-absences itself. It has negative effects on the predictive  
131 performance of HSMs and it is particularly critical for machine learning techniques, while  
132 regression techniques such as generalised linear models seem to be less affected (Barbet-  
133 Massin et al., 2012; Grimmer, Whitsed and Horta, 2020; Valavi et al., 2021). So far, class  
134 overlap has been addressed using resampling techniques more oriented to adjusting an  
135 unbalanced number of classes in the response variable (i.e., the ‘up-’ or ‘down-sampling’  
136 approach; Valavi et al., 2021), irrespective of the technique used to obtain pseudo-  
137 absences.

138       As far as we know, there are no approaches for sampling pseudo-absences that  
139 seek to mitigate both sample location bias and class overlap. Here, we present an  
140 alternative sampling strategy, which we called the ‘uniform’ approach, that builds upon  
141 existing strategies for systematically sampling the environmental space to select pseudo-  
142 absences. The novel aspect of the uniform approach is that, beyond reducing sample  
143 location bias, it also minimises class overlap by implementing a kernel-based filter that is  
144 used to delineate the portion of the environmental space where to collect pseudo-absences.  
145 To test our approach, we simulated 50 virtual species and compared the predictive  
146 performance of HSMs trained on pseudo-absences sampled using the uniform approach as  
147 well as other sampling strategies traditionally carried out within the geographical space:  
148 random (i.e., pseudo-absences randomly sampled within the geographical space) and  
149 buffer-out (i.e., pseudo-absences randomly collected outside buffers built around presence  
150 locations). To foster reproducibility, we provide an accompanying R package called USE  
151 (Uniform Sampling of the Environmental space), which bundles the R functions needed to  
152 implement the uniform approach. The package is available at  
153 <https://github.com/danddr/USE>. Finally, we provide a tutorial to explain how to apply the  
154 uniform approach to real case studies, using the European beech *Fagus sylvatica* L. as a



USE: a novel approach to uniformly sample the environmental space

155 target species.

156

## 157 2 Methods

### 158 2.1 Simulation of virtual species

159 We used virtual species (hereafter VS), a simulation tool that provides the great advantage  
160 of knowing the true generative process underlying the species geographical distribution  
161 (Meynard et al., 2019). We created the realised environmental space (*sensu* Jackson and  
162 Overpeck 2000) of 50 different virtual species using the bioclimatic variables gathered from  
163 the WorldClim database ([www.worldclim.org](http://www.worldclim.org); spatial resolution ~18.6 km at the Equator;  
164 Fick and Hijmans, 2017). We restricted the distribution of the simulated VS (and those of  
165 the bioclimatic variables) to the geographical extent spanning from -12° W to 25° E and  
166 from 36° to 60° N (approximately Western and Southern Europe) to significantly reduce the  
167 computational effort to process the entire workflow. Each VS was generated using a  
168 random set of five bioclimatic variables (out of the 19) through the function  
169 `generateRandomSp` from the R package `virtualspecies` (Leroy et al., 2016), which  
170 randomly assigns relationships between the VS and the bioclimatic variables (e.g., linear,  
171 quadratic relationships). This way, we obtained a raster layer reporting the habitat suitability  
172 index of each VS (HSI, Fig. 1a), which we then converted to a binary (i.e., presence-  
173 absence) map using the function `convertToPA`. Further details about parameters setting  
174 can be found in the R code available at [https://github.com/danddr/USE\\_paper](https://github.com/danddr/USE_paper).

### 175 2.2 Sampling of the pseudo-absences

USE: a novel approach to uniformly sample the environmental space

176 Regardless of the sampling approach and modelling technique used to calibrate the HSMS,  
177 the ratio between the number of presences and pseudo-absences in the training datasets  
178 (i.e., sample prevalence) was kept equal to 1, which means that an equal number of  
179 presences and pseudo-absences were collected. In practice, each of the VS-specific  
180 training dataset included 300 presences, which were randomly sampled within the  
181 geographical extent using the function `sampleOccurrences` from the `virtualspecies`  
182 R package. Consequently, we collected an equal number of pseudo-absences according to  
183 the three sampling strategies presented below.

#### 184 *2.2.1 Uniform approach: pseudo-absences sampled within the environmental space*

185 For each VS (i.e., iteration), we built a 2-dimensional environmental space by keeping the  
186 first two axes of a principal component analysis (PCA) performed on the correlation matrix  
187 of the five randomly selected bioclimatic variables used to generate the realised  
188 environment (Fig. 1b). Each time, we checked that the first two principal component axes  
189 accounted for at least 70% of the total bioclimatic variability. Then, we uniformly sampled  
190 pseudo-absences in the environmental space using the `uniformSampling` function. In  
191 short, each pseudo-absence is associated with a geographical location (i.e., a pixel of the  
192 environmental layers), which is in turn characterised by the set of environmental conditions  
193 encountered at that location. Such a combination of environmental conditions determines  
194 the position of the pseudo-absence within the environmental space. A pseudo-absence can  
195 thus be defined as the projection of a geographical location onto the environmental space  
196 generated through the PCA (i.e., a PC-score). Below, we present a step-by-step description  
197 of the uniform sampling performed by the function `paSampling`, which internally calls  
198 `uniformSampling` (both functions are included in the USE R package):

USE: a novel approach to uniformly sample the environmental space

199 1. First, kernel density estimation (a statistical technique used to estimate the underlying  
200 probability distribution of a set of data points by smoothing them with a kernel  
201 function; Scott, 1992) is used to calculate the probability density function of the  
202 presence data within the 2-dimensional environmental space. Similar uses of kernel  
203 density estimation have become popular in recent years, especially due to their  
204 increasing use in trait-based ecology to compute probabilistic hypervolumes and trait  
205 probability densities (Mammola and Cardoso, 2020 and reference therein). The PC-  
206 scores associated with a probability threshold equal to or greater than 0.75 (i.e., the  
207 default threshold value used in the `paSampling` function) are likely to bear  
208 environmental conditions associated with presence locations. Thus, we selected these  
209 presence locations and we generated the convex hull delimiting the portion of the  
210 environmental space mostly associated with this set of presence points within the  
211 environmental space (Fig. 1c). The kernel bandwidth (i.e., the width of the kernel  
212 density function that defines its shape) can be either defined by the user or  
213 automatically estimated by the function `paSampling`. In the latter case, the function  
214 uses a bandwidth selector by internally calling the function `Hpi` of the R package `ks`  
215 (Duong, 2021).

216 2. The portion of the environmental space defined by the above-mentioned convex hull is  
217 removed from the whole environmental space. Then, a sampling grid was generated  
218 from a pre-selected resolution (e.g.,  $10 \times 10$  cells) and overlaid on the 2-dimensional  
219 environmental space (Fig. 1d). The optimal resolution of the sampling grid within the  
220 environmental space can be determined using the function `optImRes` from the USE  
221 package. This function operates as follows:

222 - Within each cell of the sampling grid, the average (squared) Euclidean distance

USE: a novel approach to uniformly sample the environmental space

223 between the pseudo-absences (PC-scores) in the cell and the centroid of their  
224 convex hull is computed;

225 - Once this metric is computed across all cells of the sampling grid, the average mean  
226 value is computed across all cells (hereafter, grid average);

227 - The procedure above is separately repeated on different sampling grids of  
228 increasing resolution (i.e., increasing number of cells);

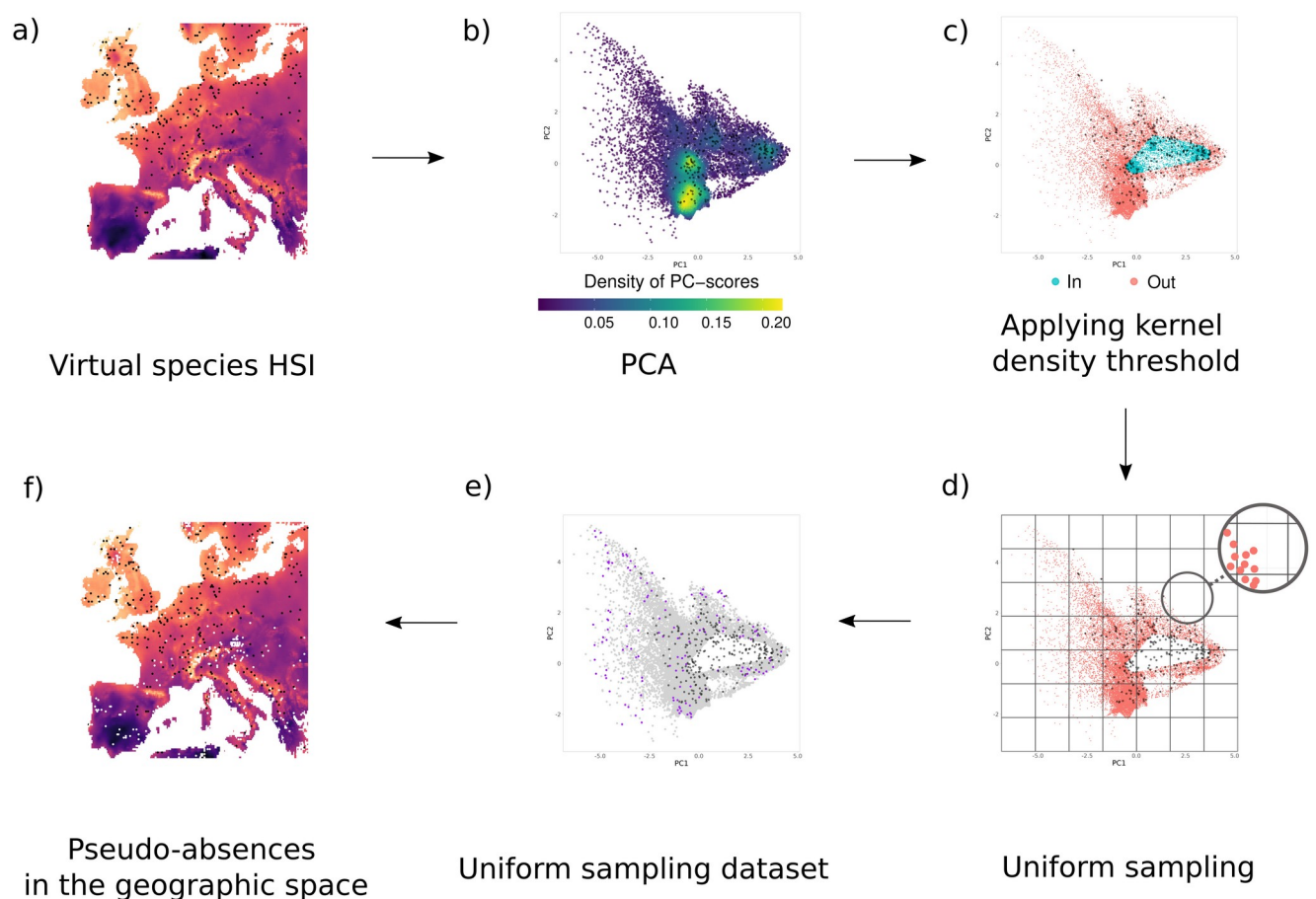
229 - The resulting set of grid averages (one per resolution) are used as a measure of the  
230 aggregation among pseudo-absences within the cells of the sampling grids. This  
231 value is compared across resolutions and the best grid is chosen as the one  
232 providing the best trade-off between resolution and average distance among points  
233 within cells (i.e., the resolution that allows uniformly sampling the environmental  
234 space without overfitting it). More specifically, the best grid is the one whose  
235 resolution is just below that which would not allow the average distance among  
236 pseudo-absences to be reduced by more than 10% (other values can be set by the  
237 user).

238 3. Once the optimal resolution is set, the sampling grid is sequentially scanned (i.e., cell  
239 by cell) by the `uniformSampling` function called via the `paSampling` function and,  
240 from each grid cell, a given number of pseudo-absences is randomly collected. At this  
241 stage, the pseudo-absences associated with environmental conditions too close to  
242 those of the presence locations are already excluded (see step 1). Note that the  
243 pseudo-absences are randomly selected within the area of each cell of the sampling  
244 grid, and not at the centroid nor at the nodes.

245 The total number of pseudo-absences sampled within each cell of the sampling grid can be  
246 set by the user (using the argument `n.tr`, default `n.tr = 5`), who can also indicate a

USE: a novel approach to uniformly sample the environmental space

247 desired sample prevalence. If the sample prevalence is not specified, fewer pseudo-  
248 absences are likely to be eventually sampled than expected (i.e.,  $n \cdot tr \times$  number of cells).  
249 This happens because (i) no pseudo-absence points are collected in empty cells, and (ii)  
250 less pseudo-absence points than  $n \cdot tr$  are available within the cells at the boundary of the  
251 environmental space (see zooming window in Figure 1d). Similarly, no pseudo-absences  
252 are collected within the core area of the presences (excluded in step 1). If a sample  
253 prevalence is set by the user, the sampling grid is surveyed until the chosen sample  
254 prevalence is reached by the algorithm.



255 **Figure 1:** Flowchart representing the step-by-step procedure for implementing the uniform

USE: a novel approach to uniformly sample the environmental space

256 approach: a) habitat suitability index (HSI) of the  $i$ -th virtual species (VS; lighter colours  
257 indicate higher habitat suitability and black dots represent presence points in the  
258 geographical space); b) PCA performed on the environmental variables in the study region  
259 (lighter colours indicate high PC-scores densities and black dots represent the presence  
260 points within the environmental space); c) application of the kernel-based filter, which splits  
261 the environmental space into two sub-spaces associated with either the environmental  
262 conditions more suitable for the species (in blue) or those associated with less/not suitable  
263 environmental conditions (in red; with black dots still depicting presence points); d) pseudo-  
264 absences are uniformly sampled across a sampling grid of a chosen resolution overlaid to  
265 the 2-dimensional environmental space. Specifically, pseudo-absences are sampled within  
266 each cell of the 2-d grid. The inset map shows an example of a grid cell at the boundary of  
267 the environmental space (i.e., a grid cell containing low density of pseudo-absences), black  
268 dots represent presence points; e) the purple dots represent the pool of randomly selected  
269 pseudo-absences after running the uniform sampling approach; f) the white dots represent  
270 the selected set of pseudo-absences after running the uniform sampling approach, but  
271 displayed in the geographical space this time, black dots still represent presence points  
272 from the focal virtual species.

### 273 *2.2.2 Pseudo-absences sampled within the geographical extent*

274 The sampling of pseudo-absences within the geographical extent was conducted using the  
275 random and buffer-out approaches. For the random approach (Barbet-Massin et al. 2012;  
276 Iturbide et al., 2015; Støa et al., 2019), we simply generated 300 random pseudo-absences  
277 across the studied geographical extent. For the buffer-out approach (Bedia et al., 2013), we  
278 created a buffer of 50 km radius around each presence location, and we then randomly  
279 sampled pseudo-absences outside the presence-specific buffers, but within the convex hull

USE: a novel approach to uniformly sample the environmental space

280 of the species geographical distribution (i.e., the convex hull that connects the outer  
281 presences of the species and thus delimits the range actually covered by the species in the  
282 geographical space).

### 283 2.3 Habitat suitability models

284 For each of the 50 VS and for each of the three sampling strategies (i.e., uniform, random,  
285 buffer-out), we built a specific dataset combining the presence records with the pseudo-  
286 absences sampled within the environmental and the geographical space. First, we modelled  
287 the presence and pseudo-absences data as a function of the same five bioclimatic variables  
288 used to generate each of the 50 VS. To this aim, we randomly partitioned each dataset  
289 (specific for a sampling strategy) into 5 replicates of both training (70% observations) and  
290 testing (30%) sets, which we used to calibrate and validate, respectively and for each  
291 replicate, five modelling algorithms: (i) binomial generalised linear models with 'logit' link  
292 (GLMs); (ii) generalised additive models (GAMs); (iii) random forests (RFs); (iv) boosted  
293 regression trees (BRTs); and (v) MaxEnt. In total, we fitted 3,750 HSMs (50 VS species  $\times$  3  
294 different sets of pseudo-absences  $\times$  5 modelling algorithms  $\times$  5 replicates of 70-30%  
295 partitions). To fit the HSMs, we used the R package `sdm` (Naimi and Araújo, 2016).  
296 Although we acknowledge the importance of fine-tuning HSMs (Fourcade, 2021), we kept  
297 model settings at their default value since it would have been unfeasible to individually  
298 parametrise each algorithm for all 50 VS and sampling strategies. A detailed representation  
299 of the workflow of the analyses is shown in Fig. 2. Furthermore, we acknowledge that our  
300 use of MaxEnt did not conform with the general recommendations for its adequate  
301 implementation (e.g., using 10,000 background points; Cobos et al., 2019; Kass et al.,  
302 2021). Nonetheless, we included it in the comparison of models' performance due to its  
303 wide usage within the HSMs community.

USE: a novel approach to uniformly sample the environmental space

## 304 2.4 Comparison among sampling strategies

### 305 *2.4.1 Predictive performance comparison*

306 After fitting HSMs for all the 50 VS, we compared the predictive performance  
307 associated with each combination of sampling approaches and modelling techniques by  
308 computing the following metrics: (i) the area under the receiver operating characteristic  
309 curve (AUC); (ii) the continuous Boyce index (CBI); (iii) the sensitivity; (iv) the specificity; (v)  
310 the true skill statistics (TSS); and (vi) the root mean squared error (RMSE). The RMSE was  
311 computed by comparing the true (i.e., simulated) habitat suitability of the focal VS against  
312 the one predicted by each combination of modelling and sampling approach. A detailed  
313 description of the above-mentioned modelling techniques and validation metrics can be  
314 found in Guisan et al. (2017). To compare the predictive performance of the HSMs fitted  
315 under different combinations of sampling strategy and modelling technique, we visually  
316 assessed the results of the 50 VS simulations using violin plots reporting the distribution of  
317 the values of the predictive performance metrics listed above. Furthermore, we tested for  
318 statistical differences between the three sampling strategies for each predictive accuracy  
319 metric using the Kruskal-Wallis test, followed by one-tailed Dunn's post hoc rank sum  
320 comparisons using the `dunn.test` R package (Dinno, 2017) (p-values for multiple  
321 comparisons adjusted using Holm correction).

322

### 323 *2.4.2 Sample location bias and class overlap*

324 To assess the intensity of sample location bias associated with the different sampling  
325 strategies, we extracted the pseudo-absences of a single VS and map their aggregation  
326 within the environmental space using bivariate density plots. The aim was to identify which,



USE: a novel approach to uniformly sample the environmental space

327 among the three sampling strategies, was more subject to oversampling particular  
328 environmental conditions within the geographical space. In principle, the sampling  
329 strategies more affected by sample location bias would exhibit a clear aggregation of  
330 pseudo-absences within the environmental space. We visually assessed the areas of the  
331 environmental space sampled by the different sampling strategies using the function  
332 `geom_density_2d` of the `ggplot2` R package (Wickham, 2016). This function performs a  
333 2D kernel density estimation using the `kde2d` function of the `MASS` R package (Venables  
334 and Ripley, 2002) and displays the results with contours. In addition, for 10 new VS, we  
335 calculated the total range (i.e., max PC-score – min PC-score) of the two principal  
336 component axes associated with the pseudo-absences collected through the different  
337 sampling strategies. We then derived the 95% confidence interval of the total range through  
338 a nonparametric bootstrap ( $n = 2,000$ ) using the function `smean.cl.boot` from the `Hmisc`  
339 R package for each principal component axis and sampling strategies. We tested for  
340 statistical differences for each principal component axis among sampling strategies using  
341 the Kruskal-Wallis test followed by two-tailed Dunn's post hoc rank sum comparisons with  
342 Holm's correction. To assess the effectiveness of the uniform approach for mitigating class  
343 overlap, we simulated 10 new VS, sampled their presences and pseudo-absences using the  
344 three sampling strategies and mapped the position of the presence and pseudo-absence  
345 points within the environmental space following the procedure explained in section 2.2.1  
346 and Figure 1a,b. Then, we computed the Gaussian hypervolume of the presences and  
347 pseudo-absences using the `hypervolumes` R package (Blonder et al., 2014; 2022), and  
348 calculated the overlap between them. Statistically significant differences in the degree of  
349 overlap were tested using one-way ANOVA and Tukey HSD test.

USE: a novel approach to uniformly sample the environmental space

## 350 2.5 Sensitivity analyses

351 In our analytical framework, we kept fixed the value of the following parameters: sample  
352 prevalence, the size of the buffer for the buffer-out approach, and the number of bioclimatic  
353 variables used as predictors to fit the HSMs for the VS. To test the potential effect on our  
354 results of varying these parameters, we conducted the following sensitivity analyses:

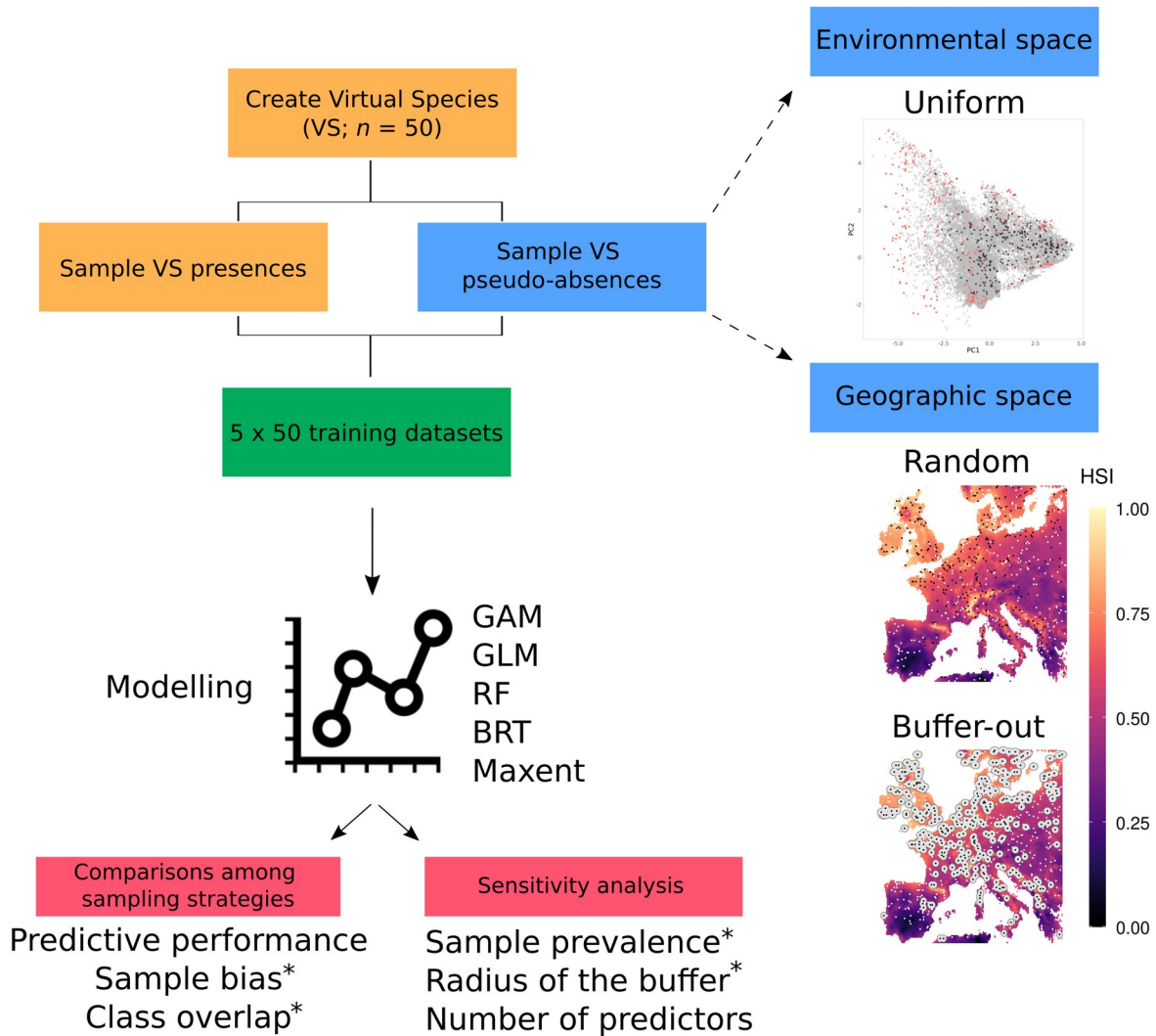
- 355 • To test the effect of changing sample prevalence on the predictive performance of the  
356 different sampling strategies, we repeated the entire workflow on 10 VS using two additional  
357 prevalence values, namely 0.5 and 0.1. Specifically, for each VS, we generated two  
358 additional training datasets with 300 presences, but we combined them with 600 and 3,000  
359 pseudo-absences to achieve sample prevalence of 0.5 and 0.1 respectively.
- 360 • To test the effect of the size of the buffer on the predictive performance of the buffer-out  
361 approach, we repeated the entire workflow on 10 VS considering the following buffer radius  
362 lengths: 50, 100 and 200 km.
- 363 • To test how using a different number of bioclimatic variables would affect the predictive  
364 performance of the sampling strategies, we repeated the entire workflow on 50 VS using all  
365 19 bioclimatic variables to both define the environmental space to generate the VS and as  
366 predictors to fit the related HSMs.

## 367 2.6 Real-case study

368 To illustrate how to apply the uniform approach with the USE R package, we modelled the  
369 realised distribution of *Fagus sylvatica* in Italy, France and Spain. We chose *F. sylvatica* as  
370 a target species because its distribution and biogeographic history is well-known across  
371 Europe (Magri et al., 2006; Poli et al., 2022). The whole analysis of *F. sylvatica* is described

USE: a novel approach to uniformly sample the environmental space

372 in S5, and the R code to replicate it can be found at: [https://github.com/danddr/USE\\_paper](https://github.com/danddr/USE_paper).



373 **Figure 2** Overall workflow of the analysis described in the Methods section. The '\*' is  
 374 associated with analyses (i.e., sample bias, class overlap, sample prevalence, radius of the  
 375 buffer) performed on  $n = 10$  virtual species (VS).

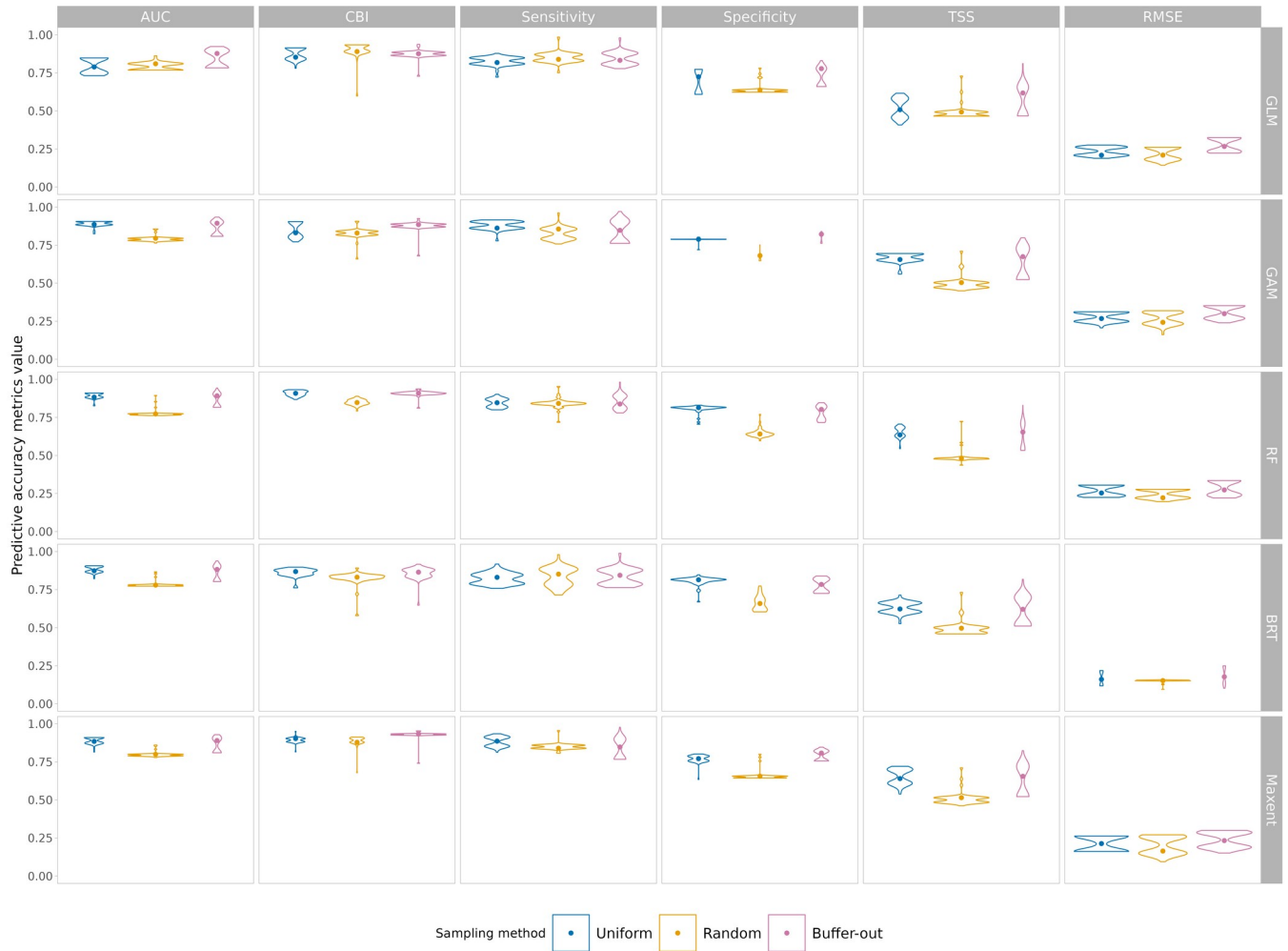
USE: a novel approach to uniformly sample the environmental space

## 376 3 Results

### 377 3.1 Comparison of the predictive performance associated with geographical vs 378 environmental sampling

379 Overall, the uniform approach performed equal to or better than the geographical  
380 approaches in terms of out-of-sample prediction (Fig. 3). Pairwise comparisons between the  
381 predictive accuracy performance of the uniform approach against the random and buffer-out  
382 approaches showed statistically significant differences in 73% and 47% of the  
383 combinations, respectively. However, these differences were algorithm- and metric-  
384 dependent and did not point to an overall higher predictive performance of the uniform  
385 approach (Fig. 3, Tab. S1, Fig. S1.1). The pattern of the differences among predictive  
386 performance metrics was consistent among prevalence values (Fig. S2.1-2.2) and number  
387 of bioclimatic variables used in the models (Fig. S3). Increasing the buffer radius length  
388 (Fig. S4), resulted in higher predictive performance of the buffer-out approach for some  
389 metrics (AUC, TSS, specificity), while for CBI, sensitivity and RMSE results remained  
390 comparable with those presented in Fig. 3.

USE: a novel approach to uniformly sample the environmental space



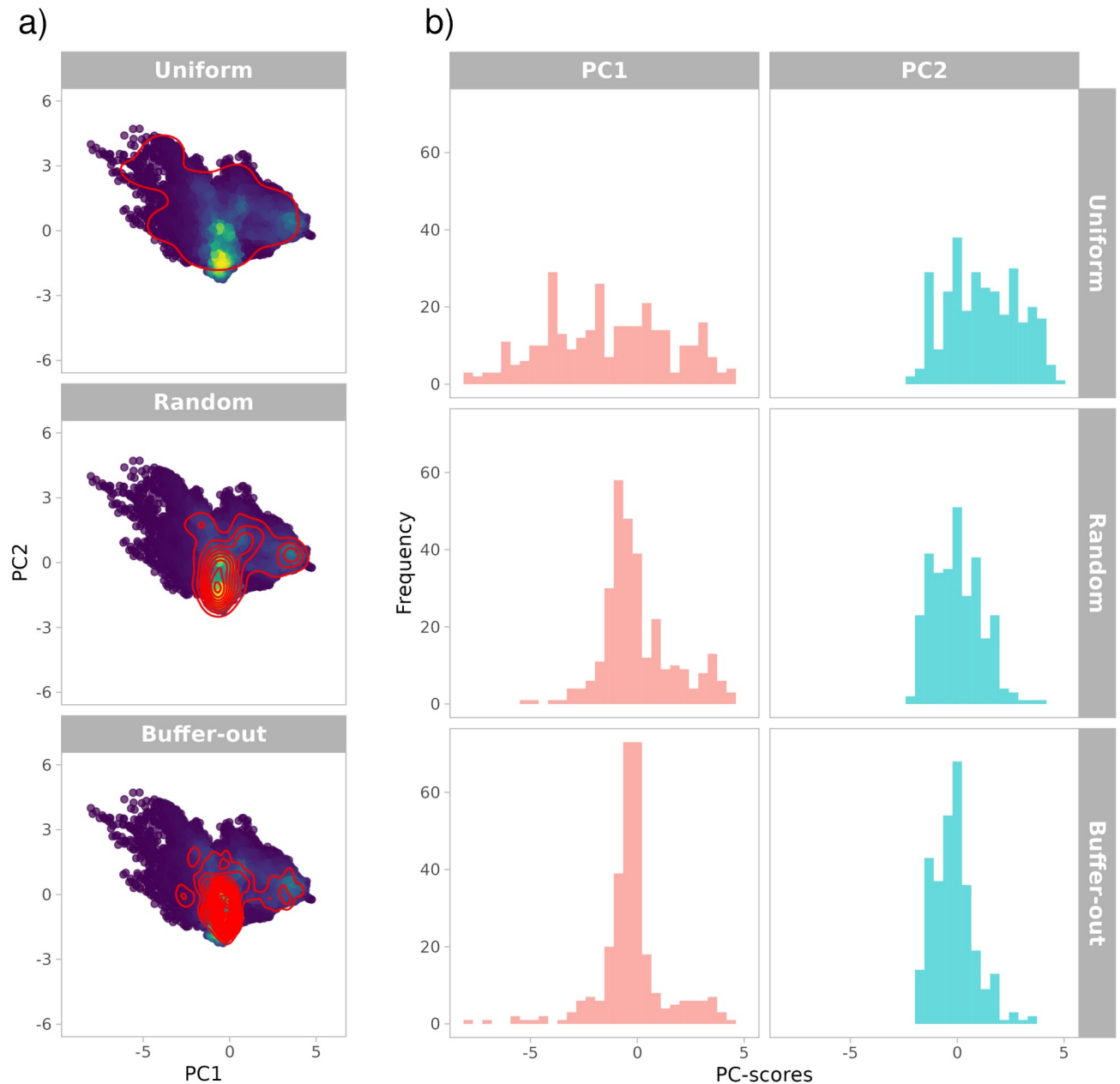
391 **Figure 3:** Violin plots reporting the distribution of the values of the metrics of predictive  
 392 performance for the habitat suitability models (HSMs) of the 50 virtual species (VS), as modelled  
 393 using 5 randomly selected bioclimatic predictors and setting sample prevalence equal to 1 (i.e.,  
 394 same number of presences and pseudo-absences). Dots represent median values of the metrics of  
 395 predictive accuracy. Columns indicate the different performance metrics, while rows are associated  
 396 with the modelling techniques used to fit the HSMs. Higher values in all metrics but RMSE reflect  
 397 higher predictive performance. AUC = area under the curve; CBI = continuous Boyce index, TSS =  
 398 true skill statistic; RMSE = root mean squared error; GLM = generalised linear model; GAM =  
 399 generalised additive model; RF = random forest; BRT = boosted regression trees.

USE: a novel approach to uniformly sample the environmental space

### 400 3.2 Effect of sample location bias and class overlap

401 The bivariate density plots of the pseudo-absences sampled within the environmental and  
402 geographical space highlighted that the uniform approach had the widest and most  
403 homogeneous coverage of environmental conditions throughout the environmental space  
404 (Fig. 4, see Figure S1.2 for a more detailed representation of the density of pseudo-  
405 absences sampled within the environmental space when running the uniform approach;  
406 Fig.S1.3). In contrast, the random and buffer-out approaches appeared to be prone to  
407 sample location bias, with peaks of high density of pseudo-absences occurring in specific  
408 areas of the environmental space, i.e., those associated with the most frequent habitat  
409 conditions encountered within the geographical space, and a narrow mean range of PC-  
410 scores sampled along both principal component axes compared to the uniform approach  
411 (Fig. 4, Fig. S1.3; Kruskal-Wallis test for PC1:  $\chi^2= 21.54$ ,  $df = 2$ ,  $p\text{-value} < 0.001$ ; Kruskal-  
412 Wallis test for PC2:  $\chi^2= 14.91$ ,  $df = 2$ ,  $p\text{-value} < 0.001$ ).

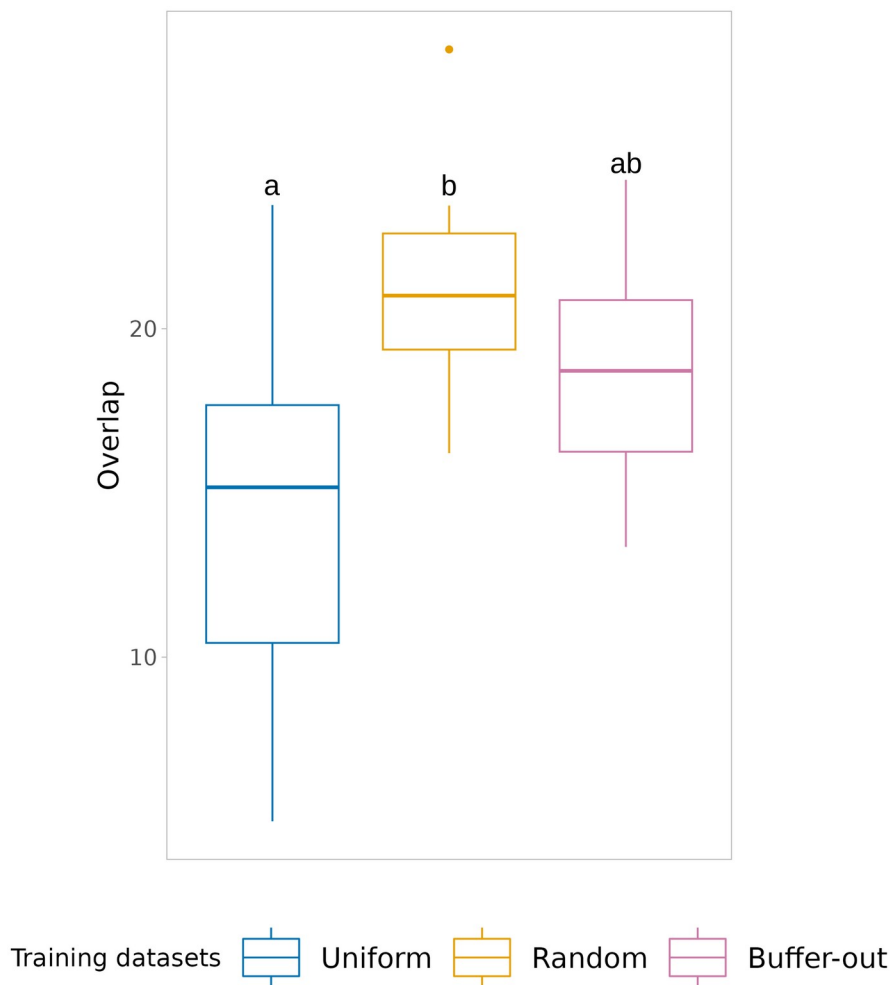
413 Regarding class overlap, we detected a statistically significant difference in the overlap  
414 between the portions of the environmental space occupied by presences and pseudo-  
415 absences sampled through different approaches (one-way ANOVA  $F(2, 27) = 5.83$ ,  $p\text{-value}$   
416  $= 0.008$ ). Specifically, the uniform approach exhibited the lowest overlap in comparison to  
417 the other sampling strategies (Fig. 5). The post hoc Tukey HSD test showed that the  
418 uniform approach exhibited a significantly lower overlap than the random sampling ( $p <$   
419  $0.001$ ), whereas the uniform- buffer-out and buffer-out-random comparisons did not show  
420 significant differences ( $p = 0.09$ ,  $p=0.47$ ).



421 **Figure 4:** a) Bivariate plots showing the environmental space generated by a principal  
 422 component analysis carried out on 5 bioclimatic variables. Red lines represent the density  
 423 of pseudo-absences for an individual virtual species, as sampled by the random and buffer-  
 424 out approaches within the geographical space, and by the uniform approach within the  
 425 environmental space. A more detailed representation of the density of pseudo-absences  
 426 sampled by the uniform approach is reported in Figure S1.2. b) Histograms showing the

USE: a novel approach to uniformly sample the environmental space

427 frequency distribution of the first two principal components (columns) associated with the  
428 different sampling strategies (rows).



429 **Figure 5:** Box plots showing the overlap between environmental spaces generated by  
430 presences and pseudo-absences of the virtual species. Letters denote significant  
431 differences using Tukey HSD test. Colours are associated with the three sampling  
432 strategies used to generate the pseudo-absences (uniform in blue, random in yellow and  
433 buffer-out in pink).

434



USE: a novel approach to uniformly sample the environmental space

## 435 4 Discussion

436 In this study, we proposed the uniform approach as an alternative strategy to sample  
437 pseudo-absences within the environmental space. In contrast to existing techniques, our  
438 approach systematically samples pseudo-absences from portions of the environmental  
439 space excluding the conditions that are likely to be suitable for the species to establish. As a  
440 result, the uniform approach reduces the chance of including false-absences in the training  
441 dataset. From a more theoretical perspective, data collected after the application of the  
442 kernel-based filter are much closer to the concept of pseudo-absences than those obtained  
443 through traditional, geographical sampling approaches. Our findings show that the uniform  
444 approach represents a valid strategy for gathering pseudo-absences, resulting in out-of-  
445 sample predictive accuracy comparable to the sampling strategies implemented within the  
446 geographical space. In addition, the uniform sampling significantly reduces sample location  
447 bias and class overlap, which is critical to obtain ecologically meaningful pseudo-absences.  
448 Importantly, the uniform approach is flexible, as it allows the user to set parameters (e.g.,  
449 kernel bandwidth, sample prevalence, sampling grid resolution) that control how pseudo-  
450 absences are sampled within the environmental space. Such flexibility is particularly  
451 valuable to mimic different ecological processes that are easier to capture within the  
452 environmental space than within the geographical space (e.g., source-sink dynamics). In all  
453 cases, by generating informative pseudo-absences, the uniform approach allows satisfying  
454 one of the most critical assumptions underpinning habitat suitability modelling: the need for  
455 adequate species distribution attributes (i.e., pseudo-absence data here) to model the  
456 species-environment relationship (Guisan et al., 2017).

USE: a novel approach to uniformly sample the environmental space

#### 457 4.1 Effect of the sampling approaches on models' predictive performance

458 Results of the VS' simulations showed that the uniform approach performed well in terms of  
459 out-of-sample prediction regardless of the modelling technique, metric of predictive  
460 performance, and sample prevalence used. All HSMs calibrated on pseudo-absences  
461 sampled with the uniform approach consistently showed high predictive performance,  
462 especially for the metrics related to the capacity of a model to correctly predict presences  
463 (i.e., sensitivity and CBI). Concerning the metrics associated with the models' ability to  
464 predict absences (e.g., specificity), the uniform sampling showed values comparable to the  
465 other strategies. This suggests that the uniform approach reduces omission error without  
466 necessarily increasing commission error. This is coherent with Fei and Yu (2016), who  
467 reported an increase in overall model predictive performance when pseudo-absences were  
468 systematically collected within the environmental space.

469 In this sense, results for the CBI, which is currently the go-to accuracy metric for validating  
470 HSMs fitted on pseudo-absences (or background points), and for the RMSE were  
471 particularly encouraging since the uniform approach scored, together with the buffer-out  
472 approach, the highest CBI values and lowest RMSE values across all modelling techniques.  
473 The high predictive performance associated with the uniform approach can be attributed to  
474 its two main underlying properties: the systematic sampling of the environmental space and  
475 the kernel-based filter on the presence data.

476 Notwithstanding the positive results obtained in terms of predictive performance, we argue  
477 that a comparison of metrics of model predictive accuracy may not be the best means for  
478 evaluating the adequacy of different sampling strategies carried out within the  
479 environmental rather than the geographical space. Indeed, previous studies showed that

USE: a novel approach to uniformly sample the environmental space

480 these metrics are affected by several factors, including sample prevalence (Guisan et al.,  
481 2017; Leroy et al., 2018; Marchetto et al., 2023), sample bias (Dubos et al., 2022, Rocchini  
482 et al., 2023) or the spatial extent of the study area (Lobo et al., 2008). Moreover, AUC and  
483 TSS tend to score high even in case of poor models calibrated on data exhibiting strong  
484 sample location bias (Fourcade et al., 2018, Jiménez-Valverde, 2021). Assessing HSMS  
485 predictive performance using a set of different predictive accuracy metrics might help the  
486 user to critically evaluate the outputs of the models.

## 487 4.2 Effect of the uniform sampling on sample location bias and class 488 overlap

489 The uniform approach proved to significantly reduce sample location bias, since pseudo-  
490 absences were homogeneously scattered across the bivariate density plot of the two  
491 principal component axes (Fig. 4a,b, Fig. S1.2 in Supplementary Materials) and collected a  
492 wider range of PC-scores compared to the random and buffer-out approaches (Fig. S1.3).  
493 On the contrary, the two sampling approaches carried out within the geographical space  
494 exhibited prominent peaks of density of pseudo-absences in correspondence with the most  
495 frequently encountered environmental conditions within the geographical space, resulting in  
496 a narrower mean of PC-scores. As a consequence, the random and buffer-out approaches  
497 may provide sub-optimal pseudo-absences for modelling the species-environment  
498 relationship (Thuiller et al. 2004; Austin 2007). This aspect gets increasingly relevant as  
499 environmental conditions are more heterogeneously distributed across the geographical  
500 space (e.g., in mountain regions with high topographic heterogeneity). Therefore, HSMS  
501 calibrated on training datasets adequately representing environmental variability rather than  
502 wide geographical coverage represent a crucial step to better capture and discriminate  
503 species niche breadth (Tessarolo et al., 2014, 2021; Varela et al., 2014; Bazzichetto et al.,

USE: a novel approach to uniformly sample the environmental space

504 2022; Perret and Sax 2022).

505 The uniform approach proved to also significantly reduce class overlap. The thresh  
506 argument passed to the `paSampling` function controls the portion of the environmental  
507 space associated with the species presence, thus inherently limiting class overlap by the  
508 exclusion of environmental conditions suitable to the species (see Fig. 1c, Fig. 5 and Fig.  
509 S1.4). This results in a set of pseudo-absences theoretically much closer to the species'  
510 true absences. Given that presence points are unevenly distributed within the environmental  
511 space, different kernel thresholds might also be used to handle the sampling of pseudo-  
512 absences under particular scenarios. As an example, setting a low kernel threshold would  
513 allow excluding accidental presences from unsuitable locations (e.g., 'sink populations')  
514 from the training dataset, while potentially including observations from these areas as  
515 pseudo-absences. Unfortunately, there is no *a priori* choice about the value of the threshold  
516 without having preliminary information on species' ecology, the study area and the goal of  
517 the research. For this reason, we provided the `thresh.inspect` function, which produces  
518 plots depicting the entire environmental space alongside the portion that would be excluded  
519 based on a specific kernel density threshold.

## 520 4.4 Limitations and usage notes

### 521 4.4.1 Limitations

522 The first limitation of the uniform approach, which is anyway a general limitation in HSMs  
523 (e.g., Cayuela et al., 2009), is that its effectiveness depends on the amount (sample size)  
524 and quality (e.g., geographically unbiased data *sensu* Fourcade 2014) of presence data.  
525 Indeed, if few presence data are available and/or presence data are geographically biased,  
526 the kernel-based filter might not accurately delimit the area associated with suitable

USE: a novel approach to uniformly sample the environmental space

527 conditions for the species. As a consequence, the capacity of discriminating between  
528 suitable and unsuitable conditions of the uniform approach might be negatively affected.

529 A second limitation is that, although the uniform approach proved to be robust to  
530 varying sample prevalence, its effectiveness might diminish if a very large number of  
531 pseudo-absences is sampled (e.g., in case of low sample prevalence) (Fig. S2.1-2.2). Since  
532 the uniform approach samples a user-defined number of pseudo-absences within a grid  
533 overlaid to a bi-dimensional environmental space, if the number of pseudo-absences grows  
534 indefinitely, the advantage of the systematic sampling decreases. Indeed, oversampling the  
535 environmental space would generate datasets suffering from sample location bias as much  
536 as those based on the random sampling carried out within the geographical space.

537 From a more practical perspective, the uniform approach can currently operate only  
538 across 2-dimensional environmental spaces, but 3-dimensional spaces might be supported  
539 in the future.

540 Finally, although the idea behind USE and the uniform sampling approach is to provide  
541 users with an easy-to-use tool to generate more ecologically meaningful pseudo-absences,  
542 we acknowledge the existence of other techniques designed to avoid generating pseudo-  
543 absences altogether. Notable examples are point-process analyses (e.g., Isaac et al.,  
544 2020), which model the density of presence-only points per unit area, rather than the  
545 probability of presences and (pseudo-)absences. More recently, machine-learning methods  
546 based on isolation forests were also proposed, with the R package ITSDM specifically  
547 dedicated to HSMs (Song and Estes, 2023). We believe, however, that our approach  
548 provides a simpler and more intuitive way to deal with the issue of presence-only data, and  
549 thus has a lower threshold for end-users to implement in their workflow.

#### 550 *4.4.2 Usage notes*

551 We here used the uniform approach to sample bioclimatic spaces, although we stress the

USE: a novel approach to uniformly sample the environmental space

552 importance of not only using bioclimatic variables, but also information on soil, land-use as  
553 well as other relevant variables when modelling species distributions. Also, we invite  
554 potential users of the uniform sampling approach to always check that the first two axes of  
555 the principal component analysis used to generate the environmental space explains a  
556 large portion of the variance observed in the data (e.g.,  $\geq 70\%$ ). Equally important is the  
557 choice of the boundaries of the geographical extent for which the 2-dimensional space has  
558 to be generated. Indeed, to avoid the "there are no elephants in the Antarctic" paradox  
559 (Lobo et al., 2010), the spatial extent of the study area should be delineated so that it  
560 excludes geographical locations, and in turn environmental conditions, less suitable for the  
561 species (e.g., collecting pseudo-absences from Mediterranean coastal dunes when  
562 modelling the distribution of an alpine plant species). In short, the uniform approach can  
563 provide exhaustive information on where the species is likely to not occur, but it remains a  
564 responsibility of the end user to carefully verify if such information is ecologically  
565 meaningful.

566

## 567 5 Conclusion

568 In this study, we compared the predictive performance of two strategies for sampling pseudo-  
569 absences carried out within the geographical space with that of the uniform approach, which  
570 operated within the environmental space. Also, we compared geographical and environmental  
571 sampling approaches in terms of their vulnerability to sample location bias and class overlap. The  
572 uniform approach proved to have good predictive performances and to reduce sample location  
573 bias and class overlap, thereby representing a valid alternative to generate pseudo-absences for  
574 HSMs. We made the uniform approach openly available to the modellers community at  
575 <https://github.com/danddr/USE>.

USE: a novel approach to uniformly sample the environmental space

## 576 6 Declaration

- 577 • Ethics approval and consent to participate: Not applicable.
- 578 • Competing interests: No conflict of interest has been declared by the authors.
- 579 • Funding: DDR was supported by a FRS-FNRS Belgian grant, ET is supported by the  
580 Estonian Research Council grant (MOBJD1030), MB acknowledges funding from the  
581 European Union's Horizon Europe research and innovation programme under the Marie  
582 Skłodowska-Curie grant agreement No 101066324.
- 583 • Authors' contribution: MB conceived the idea of the uniform approach and wrote the  
584 related R functions, while ET and DDR integrated the kernel density-based estimation of  
585 presences and the prevalence-related settings. DDR, ET and MB performed the  
586 simulations, analysed the data and assembled the USE R package. JL, JJL, SOV, and  
587 DR critically commented on the results of the analyses and their interpretation; DDR, ET  
588 and MB led the writing of the manuscript and produced a first draft, which was further  
589 improved by all other authors.
- 590 • Acknowledgments: The authors are grateful to Prof. Joaquin Hortal, who provided  
591 constructive feedback and commented on a previous version of this manuscript. We are  
592 also grateful to the MEE's Associate Editor Prof. Luis Cayuela and the two anonymous  
593 reviewers for the very constructive comments and suggestions received during the  
594 revision process. Simulations were carried out using the facilities of the High-  
595 Performance Computing Center of the University of Tartu.

USE: a novel approach to uniformly sample the environmental space

## 596 7 Code and Data availability

597 The scripts for replicating the analyses presented in this paper are available at  
598 [https://github.com/danddr/USE\\_paper](https://github.com/danddr/USE_paper), as well as all the raw outputs of the simulations and  
599 statistical analyses (which are available as an .RDS file).

600 We provide a general tutorial to explain how to apply the USE package at  
601 [https://danddr.github.io/USE/articles/USE\\_vignette.html](https://danddr.github.io/USE/articles/USE_vignette.html). In addition, we provide a tutorial on  
602 how to apply the uniform approach based on a real species (the European beech, *Fagus*  
603 *sylvatica* L.) in S5. The R script related to the tutorial is available at  
604 [https://github.com/danddr/USE\\_paper](https://github.com/danddr/USE_paper).



USE: a novel approach to uniformly sample the environmental space

## 605 References

- 606 Albert, C. H., Yoccoz, N. G., Edwards Jr, T. C., Graham, C. H., Zimmermann, N. E., and  
607 Thuiller, W. (2010). Sampling in ecology and evolution – bridging the gap between theory  
608 and practice. *Ecography*, 33(6):1028–1037. [https://doi.org/10.1111/j.1600-](https://doi.org/10.1111/j.1600-0587.2010.06421.x)  
609 [0587.2010.06421.x](https://doi.org/10.1111/j.1600-0587.2010.06421.x)
- 610 Austin, M. (2007). Species distribution models and ecological theory: a critical assessment  
611 and some possible new approaches. *Ecological Modelling*, 200(1-2), 1-19.  
612 <https://doi.org/10.1016/j.ecolmodel.2006.07.005>
- 613 Barbet-Massin, M., Jiguet, F., Albert, C. H., and Thuiller, W. (2012). Selecting pseudo  
614 absences for species distribution models: how, where and how many? *Methods in*  
615 *Ecology and Evolution*, 3(2):327–338. <https://doi.org/10.1111/j.2041-210X.2011.00172.x>
- 616 Baker, D. J., Maclean, I. M. D., Goodall, M., and Gaston, K. J. (2022). Correlations between  
617 spatial sampling biases and environmental niches affect species distribution models.  
618 *Global Ecology and Biogeography*, 00, 1– 13. <https://doi.org/10.1111/geb.13491>
- 619 Batista, E., Lopes, A., Miranda, P., and Alves, A. (2023). Can species distribution models be  
620 used for risk assessment analyses of fungal plant pathogens? A case study with three  
621 Botryosphaeriaceae species. *European Journal of Plant Pathology*, 165(1), 41-56.  
622 <https://doi.org/10.1007/s10658-022-02587-7>
- 623 Bazzichetto, M., Lenoir, J., Da Re, D., Tordoni, E., Rocchini, D., Malavasi, M., Barták, V.  
624 and Sperandii, M. G. (2022). Sampling strategy matters to accurately estimate response  
625 curves' parameters in species distribution models. *Global Ecology and Biogeography*  
626 <https://doi.org/10.1111/geb.13725>
- 627 Bazzichetto, M., Massol, F., Carboni, M., Lenoir, J., Lembrechts, J. J., Joly, R., and Renault,  
628 D. (2021). Once upon a time in the far south: Influence of local drivers and functional  
629 traits on plant invasion in the harsh sub-Antarctic islands. *Journal of Vegetation*  
630 *Science*, 32(4), e13057. <https://doi.org/10.1111/jvs.13057>
- 631 Beck, J., Böller, M., Erhardt, A., and Schwanghart, W. (2014). Spatial bias in the GBIF  
632 database and its effect on modeling species' geographic distributions. *Ecological*  
633 *Informatics*, 19:10–15. <https://doi.org/10.1016/j.ecoinf.2013.11.002>
- 634 Bedia, J., Herrera, S., and Gutiérrez, J. M. (2013). Dangers of using global bioclimatic  
635 datasets for ecological niche modeling. limitations for future climate projections. *Global*  
636 *and Planetary Change*, 107:1–12. <https://doi.org/10.1016/j.gloplacha.2013.04.005>
- 637 Blonder, B., Lamanna, C., Violle, C., and Enquist, B. J. (2014). The n-dimensional  
638 hypervolume. *Global Ecology and Biogeography*, 23(5), 595-609.  
639 <https://doi.org/10.1111/geb.12146>
- 640 Blonder B, Morrow CB, Harris DJ, Brown S, Butruille G, Laini A, Chen D (2022).  
641 hypervolume: High Dimensional Geometry, Set Operations, Projection, and Inference  
642 Using Kernel Density Estimation, Support Vector Machines, and Convex Hulls. R  
643 package version 3.0.4, <https://CRAN.R-project.org/package=hypervolume>.

USE: a novel approach to uniformly sample the environmental space

- 644 Booth, T. H., Nix, H. A., Busby, J. R., and Hutchinson, M. F. (2014). Bioclim: the first  
645 species distribution modelling package, its early applications and relevance to most  
646 current maxent studies. *Diversity and Distributions*, 20(1):1–9.  
647 <https://doi.org/10.1111/ddi.12144>
- 648 Cayuela, L., Golicher, D. J., Newton, A. C., Kolb, M., De Albuquerque, F. S., Arets, E. J. M. M.,  
649 Alkemade, J. R. M. and Pérez, A. M. (2009). Species distribution modeling in the tropics: problems,  
650 potentialities, and the role of biological data for effective species conservation. *Tropical*  
651 *Conservation Science*, 2(3), 319-352. <https://doi.org/10.1177/194008290900200304>
- 652 Cobos, M. E., Peterson, A. T., Barve, N., and Osorio-Olvera, L. (2019). kuenm: an R package for  
653 detailed development of ecological niche models using Maxent. *PeerJ*, 7, e6281.  
654 <https://doi.org/10.7717/peerj.6281>
- 655 Da Re, D., Tordoni, E., De Pascalis, F., Negrín-Pérez, Z., Fernández-Palacios, J. M.,  
656 Arévalo, J. R., ... and Bacaro, G. (2020). Invasive fountain grass (*Pennisetum setaceum*  
657 (Forssk.) Chiov.) increases its potential area of distribution in Tenerife island under future  
658 climatic scenarios. *Plant Ecology*, 221(10), 867-882. [https://doi.org/10.1007/s11258-020-](https://doi.org/10.1007/s11258-020-01046-9)  
659 [01046-9](https://doi.org/10.1007/s11258-020-01046-9)
- 660 Dinno, A. (2017). *dunn.test: Dunn's Test of Multiple Comparisons Using Rank Sums*. R  
661 package version 1.3.5, <https://CRAN.R-project.org/package=dunn.test>.
- 662 Dubos, N., Préau, C., Lenormand, M., Papuga, G., Monsarrat, S., Denelle, P., ... and  
663 Luque, S. (2022). Assessing the effect of sample bias correction in species distribution  
664 models. *Ecological Indicators*, 145, 109487. <https://doi.org/10.1016/j.ecolind.2022.109487>
- 665 Duffy, G. A., Coetsee, B. W., Latombe, G., Akerman, A. H., McGeoch, M. A., and Chown,  
666 S. L. (2017). Barriers to globally invasive species are weakening across the  
667 Antarctic. *Diversity and Distributions*, 23(9), 982-996. <https://doi.org/10.1111/ddi.12593>
- 668 Duong, T. (2021). *ks: Kernel Smoothing*. R package version 1.13.3.  
669 <https://cran.r-project.org/web/packages/ks/index.html>
- 670 Fei, S. and Yu, F. (2016). Quality of presence data determines species distribution model  
671 performance: a novel index to evaluate data quality. *Landscape Ecology*, 31(1):31–42.  
672 <https://doi.org/10.1007/s10980-015-0272-7>
- 673 Fick, S. E. and Hijmans, R. J. (2017). WorldClim 2: new 1-km spatial resolution climate  
674 surfaces for global land areas. *International Journal of Climatology*, 37(12):4302–4315.  
675 <https://doi.org/10.1002/joc.5086>
- 676 Fourcade, Y. (2021). Fine-tuning niche models matters in invasion ecology. A lesson from  
677 the land planarian *Obama nungara*. *Ecological Modelling*, 457, 109686.  
678 <https://doi.org/10.1016/j.ecolmodel.2021.109686>
- 679 Fourcade, Y., Besnard, A. G., and Secondi, J. (2018). Paintings predict the distribution of  
680 species, or the challenge of selecting environmental predictors and evaluation statistics.  
681 *Global Ecology and Biogeography*, 27(2):245–256. <https://doi.org/10.1111/geb.12684>
- 682 Fourcade, Y., Engler, J. O., Rödder, D., and Secondi, J. (2014). Mapping species  
683 distributions with MAXENT using a geographically biased sample of presence data: a

USE: a novel approach to uniformly sample the environmental space

- 684 performance assessment of methods for correcting sampling bias. *PloS ONE*, 9(5),  
685 e97122. <https://doi.org/10.1371/journal.pone.0097122>
- 686 Grimmer, L., Whitted, R., and Horta, A. (2020). Presence-only species distribution models  
687 are sensitive to sample prevalence: Evaluating models using spatial prediction stability  
688 and accuracy metrics. *Ecological Modelling*, 431, 109194.  
689 <https://doi.org/10.1016/j.ecolmodel.2020.109194>
- 690 Guisan, A., Thuiller, W., and Zimmermann, N. E. (2017). *Habitat suitability and distribution*  
691 *models: with applications in R*. Cambridge University Press.
- 692 Hallgren, W., Santana, F., Low-Choy, S., Zhao, Y., and Mackey, B. (2019). Species  
693 distribution models can be highly sensitive to algorithm configuration. *Ecological*  
694 *Modelling*, 408:108719. <https://doi.org/10.1016/j.ecolmodel.2019.108719>
- 695 Hattab, T., Garzón-López, C. X., Ewald, M., Skowronek, S., Aerts, R., Horen, H., Brasseur,  
696 B., Gallet-Moron, E., Spicher, F., Decocq, G., et al. (2017). A unified framework to model  
697 the potential and realized distributions of invasive species within the invaded range.  
698 *Diversity and Distributions*, 23(7):806–819. <https://doi.org/10.1111/ddi.12566>
- 699 Harrell Jr F (2021). Hmisc: Harrell Miscellaneous. R package version 4.6-0,  
700 <https://CRAN.R-project.org/package=Hmisc>.
- 701 Hortal, J., Jiménez-Valverde, A., Gómez, J. F., Lobo, J. M., and Baselga, A. (2008).  
702 Historical bias in biodiversity inventories affects the observed environmental niche of the  
703 species. *Oikos*, 117(6):847–858. <https://doi.org/10.1111/ddi.12566>
- 704 Hysen, L., Nayeri, D., Cushman, S., and Wan, H. Y. (2022). Background sampling for multi-scale  
705 ensemble habitat selection modeling: Does the number of points matter?. *Ecological Informatics*,  
706 72, 101914. <https://doi.org/10.1016/j.ecoinf.2022.101914>
- 707 Isaac, N. J., Jarzyna, M. A., Keil, P., Dambly, L. I., Boersch-Supan, P. H., Browning, E., ... and  
708 O'Hara, R. B. (2020). Data integration for large-scale models of species distributions. *Trends in*  
709 *Ecology and Evolution*, 35(1), 56-67. <https://doi.org/10.1016/j.tree.2019.08.006>
- 710 Iturbide, M., Bedia, J., Herrera, S., del Hierro, O., Pinto, M., and Gutiérrez, J. M. (2015). A  
711 framework for species distribution modelling with improved pseudo-absence generation.  
712 *Ecological Modelling*, 312:166–174. <https://doi.org/10.1016/j.ecolmodel.2015.05.018>
- 713 Jackson, S. T. and Overpeck, J. T. (2000). Responses of plant populations and  
714 communities to environmental changes of the late quaternary. *Paleobiology*, 26(S4):194–  
715 220. <https://doi.org/10.1017/S0094837300026932>
- 716 Jarvie, S., and Svenning, J. C. (2018). Using species distribution modelling to determine  
717 opportunities for trophic rewilding under future scenarios of climate change. *Philosophical*  
718 *Transactions of the Royal Society B: Biological Sciences*, 373(1761), 20170446.  
719 <https://doi.org/10.1098/rstb.2017.0446>
- 720 Jiménez-Valverde, A. (2021). Prevalence affects the evaluation of discrimination capacity in  
721 presence-absence species distribution models. *Biodiversity and Conservation*, 30(5),  
722 1331-1340. <https://doi.org/10.1007/s10531-021-02144-4>

USE: a novel approach to uniformly sample the environmental space

- 723 Jiménez-Valverde, A., Lobo, J. M., and Hortal, J. (2008). Not as good as they seem: the  
724 importance of concepts in species distribution modelling. *Diversity and*  
725 *Distributions*, 14(6), 885-890. <https://doi.org/10.1111/j.1472-4642.2008.00496.x>
- 726 Jiménez-Valverde, A., Acevedo, P., Barbosa, A. M., Lobo, J. M., and Real, R. (2013).  
727 Discrimination capacity in species distribution models depends on the representativeness  
728 of the environmental domain. *Global Ecology and Biogeography*, 22(4):508–516.  
729 <https://doi.org/10.1111/geb.12007>
- 730 Kass, J. M., Muscarella, R., Galante, P. J., Bohl, C. L., Pinilla-Buitrago, G. E., Boria, R. A., Soley-  
731 Guardia, M., and Anderson, R. P. (2021). ENMeval 2.0: Redesigned for customizable and  
732 reproducible modeling of species' niches and distributions. *Methods in Ecology and*  
733 *Evolution*, 12(9), 1602-1608. <https://doi.org/10.1111/2041-210X.13628>
- 734 Leroy, B., Delsol, R., Hugueny, B., Meynard, C. N., Barhoumi, C., Barbet-Massin, M., and  
735 Bellard, C. (2018). Without quality presence–absence data, discrimination metrics such  
736 as tss can be misleading measures of model performance. *Journal of Biogeography*,  
737 45(9):1994–2002. <https://doi.org/10.1111/jbi.13402>
- 738 Leroy, B., Meynard, C. N., Bellard, C., and Courchamp, F. (2016). virtualspecies, an R  
739 package to generate virtual species distributions. *Ecography*, 39(6):599–607.  
740 <https://doi.org/10.1111/ecog.01388>
- 741 Lobo, J. M., Jiménez-Valverde, A., and Hortal, J. (2010). The uncertain nature of absences  
742 and their importance in species distribution modelling. *Ecography*, 33(1):103–114.  
743 <https://doi.org/10.1111/j.1600-0587.2009.06039.x>
- 744 Lobo, J. M., Jiménez-Valverde, A., and Real, R. (2008). AUC: a misleading measure of the  
745 performance of predictive distribution models. *Global Ecology and Biogeography*,  
746 17(2):145– 151. <https://doi.org/10.1111/j.1466-8238.2007.00358.x>
- 747 Magri, D., Vendramin, G. G., Comps, B., Dupanloup, I., Geburek, T., Gömöry, D., ... and De  
748 Beaulieu, J. L. (2006). A new scenario for the Quaternary history of European beech  
749 populations: palaeobotanical evidence and genetic consequences. *New Phytologist*,  
750 171(1), 199-221. <https://doi.org/10.1111/j.1469-8137.2006.01740.x>
- 751 Mammola, S. and Cardoso, P. (2020). Functional diversity metrics using kernel density n-  
752 dimensional hypervolumes. *Methods in Ecology and Evolution*, 11(8):986–995.  
753 <https://doi.org/10.1111/2041-210X.13424>
- 754 Marchetto, E., Da Re, D., Tordoni, E., Bazzichetto, M., Zannini, P., Celebrin, S., ... &  
755 Rocchini, D. (2023). Testing the effect of sample prevalence and sampling methods on  
756 probability-and favourability-based SDMs. *Ecological Modelling*, 477, 110248.  
757 <https://doi.org/10.1016/j.ecolmodel.2022.110248>
- 758 Meynard, C. N., Leroy, B., and Kaplan, D. M. (2019). Testing methods in species  
759 distribution modelling using virtual species: what have we learnt and what are we  
760 missing? *Ecography*, 42(12):2021–2036. <https://doi.org/10.1111/ecog.04385>
- 761 Naimi, B. and Araújo, M. B. (2016). sdm: a reproducible and extensible R platform for  
762 species distribution modelling. *Ecography*, 39(4):368–375.  
763 <https://doi.org/10.1111/ecog.01881>

USE: a novel approach to uniformly sample the environmental space

- 764 Newbold, T. (2018). Future effects of climate and land-use change on terrestrial vertebrate  
765 community diversity under different scenarios. *Proceedings of the Royal Society B*,  
766 285(1881):20180792. <https://doi.org/10.1098/rspb.2018.0792>
- 767 Perret, D. L. and Sax, D. F. (2022). Evaluating alternative study designs for optimal  
768 sampling of species' climatic niches. *Ecography*, 2022(1).  
769 <https://doi.org/10.1111/ecog.06014>
- 770 Poli et al. (2022) Coupling fossil records and traditional discrimination metrics to test how  
771 genetic information improves species distribution models of the European beech *Fagus*  
772 *sylvatica*. *European Journal of Forest Research*, 141: 253–265  
773 <https://doi.org/10.1007/s10342-021-01437-1>
- 774 Phillips, S. J., Anderson, R. P., Dudík, M., Schapire, R. E., and Blair, M. E. (2017). Opening  
775 the black box: An open-source release of Maxent. *Ecography*, 40(7):887–893.  
776 <https://doi.org/10.1111/ecog.03049>
- 777 Phillips, S. J., Dudík, M., Elith, J., Graham, C. H., Lehmann, A., Leathwick, J., and Ferrier,  
778 S. (2009). Sample selection bias and presence-only distribution models: implications for  
779 background and pseudo-absence data. *Ecological Applications*, 19(1):181–197.  
780 <https://doi.org/10.1890/07-2153.1>
- 781 Rocchini, D., Tordoni, E., Marchetto, E. *et al.* A quixotic view of spatial bias in modelling  
782 the distribution of species and their diversity. *npj biodiversity* 2, 10 (2023).  
783 <https://doi.org/10.1038/s44185-023-00014-6>
- 784 Ronquillo, C., Alves-Martins, F., Mazimpaka, V., Sobral-Souza, T., Vilela-Silva, B., Medina,  
785 N. G., and Hortal, J. (2020). Assessing spatial and temporal biases and gaps in the  
786 publicly available distributional information of Iberian mosses. *Biodiversity Data Journal*,  
787 8. <https://doi.org/10.3897/BDJ.8.e53474>
- 788 Santini, L., Benítez-López, A., Maiorano, L., Čengić, M., and Huijbregts, M. A. (2021).  
789 Assessing the reliability of species distribution projections in climate change research.  
790 *Diversity and Distributions*, 27(6):1035–1050. <https://doi.org/10.1111/ddi.13252>
- 791 Scott, D.W. (1992) *Multivariate Density Estimation: Theory, Practice, and Visualization*,  
792 John Wiley & Sons
- 793 Sillero, N. and Barbosa, A. M. (2020). Common mistakes in ecological niche models.  
794 *International Journal of Geographical Information Science*, pages 1–14.  
795 <https://doi.org/10.1080/13658816.2020.1798968>
- 796 Song, L., and Estes, L. (2023). ITSDM: Isolation forest-based presence-only species  
797 distribution modelling and explanation in R. *Methods in Ecology and Evolution*, 14(3),  
798 831-840. <https://doi.org/10.1111/2041-210X.14067>
- 799 Støa, B., Halvorsen, R., Stokland, J. N., and Gusarov, V. I. (2019). How much is enough?  
800 influence of number of presence observations on the performance of species distribution  
801 models. *Sommerfeltia*, 39(1):1–28. <https://doi.org/10.2478/som-2019-0001>
- 802 Svenning, J.-C. and Skov, F. (2004). Limited filling of the potential range in European tree  
803 species. *Ecology Letters*, 7(7):565–573. <https://doi.org/10.1111/j.1461->

USE: a novel approach to uniformly sample the environmental space

804 [0248.2004.00614.x](https://doi.org/10.1016/j.ecolind.2020.107147)

805 Tassarolo, G., Lobo, J. M., Rangel, T. F., and Hortal, J. (2021). High uncertainty in the  
806 effects of data characteristics on the performance of species distribution models.  
807 *Ecological Indicators*, 121:107147. <https://doi.org/10.1016/j.ecolind.2020.107147>

808 Tassarolo, G., Rangel, T. F., Araújo, M. B., and Hortal, J. (2014). Uncertainty associated  
809 with survey design in species distribution models. *Diversity and Distributions*,  
810 20(11):1258–1269. <https://doi.org/10.1111/ddi.12236>

811 Thuiller, W., Brotons, L., Araújo, M. B., and Lavorel, S. (2004). Effects of restricting  
812 environmental range of data to project current and future species distributions.  
813 *Ecography*, 27(2), 165– 172. <https://doi.org/10.1111/j.0906-7590.2004.03673.x>

814 Valavi, R., Elith, J., Lahoz-Monfort, J. J., and Guillera-Aroita, G. (2021). Modelling species  
815 presence-only data with random forests. *Ecography*, 44(12):1731–1742.  
816 <https://doi.org/10.1111/ecog.05615>

817 VanDerWal, J., Shoo, L. P., Graham, C., and Williams, S. E. (2009). Selecting pseudo-  
818 absence data for presence-only distribution modeling: how far should you stray from what  
819 you know?. *Ecological Modelling*, 220(4), 589-594.  
820 <https://doi.org/10.1016/j.ecolmodel.2008.11.010>

821 Varela, S., Anderson, R. P., García-Valdés, R., and Fernández-González, F. (2014).  
822 Environmental filters reduce the effects of sampling bias and improve predictions of  
823 ecological niche models. *Ecography*, 37(11):1084–1091. <https://doi.org/10.1111/j.1600-0587.2013.00441.x>

825 Venables WN, and Ripley BD (2002). *Modern Applied Statistics with S*, Fourth edition.  
826 Springer, New York. ISBN 0-387-95457-0, <https://www.stats.ox.ac.uk/pub/MASS4/>

827 Wasof, S., Lenoir, J., Aarrestad, P. A., Alsos, I. G., Armbruster, W. S., Austrheim, G., ... &  
828 Decocq, G. (2015). Disjunct populations of European vascular plant species keep the  
829 same climatic niches. *Global Ecology and Biogeography*, 24(12), 1401-1412.  
830 <https://doi.org/10.1111/geb.12375>

831 Wickham H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. ISBN  
832 978-3-319-24277-4, <https://ggplot2.tidyverse.org>.

# Supplementary Material 1

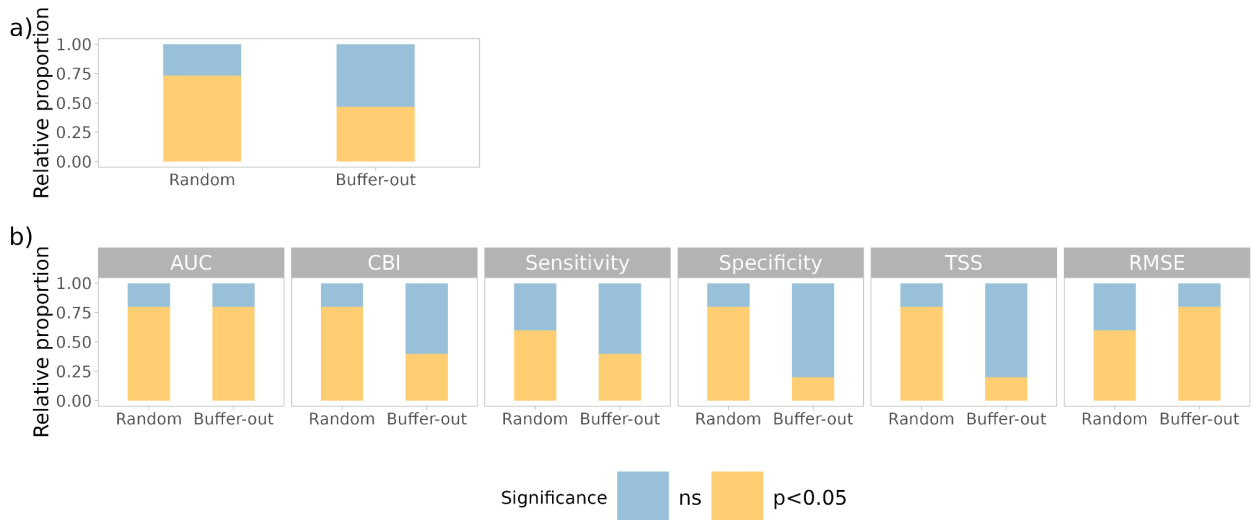
**Tab. S1:** Post-hoc multiple comparisons with Dunn’s rank sum test ( $\alpha = 0.05$ ; omnibus test was always significant with  $P < 0.05$ , data not shown). All the comparisons were performed comparing the uniform dataset against the other different sampling strategies. P-values were adjusted using Holm correction. GLM = generalised linear model; GAM = generalised additive model; RF = random forest; BRT = boosted regression trees. AUC = area under the curve; CBI = continuous Boyce index, TSS = true skill statistic; RMSE = root mean squared error. Z: test statistics. P.val: p-value (ns: not statistically significant).

Model	Metric	Comparisons	Z	P.val
BRT	AUC	Uniform - BufferOut	0.6241	ns
BRT	AUC	Uniform - Random	8.6859	p<0.05
BRT	CBI	Uniform - BufferOut	1.1292	ns
BRT	CBI	Uniform - Random	6.3851	p<0.05
BRT	RMSE	Uniform - BufferOut	-2.3726	ns
BRT	RMSE	Uniform - Random	-0.6024	ns
BRT	Sensitivity	Uniform - BufferOut	-0.9328	ns
BRT	Sensitivity	Uniform - Random	-1.375	ns
BRT	Specificity	Uniform - BufferOut	1.7994	ns
BRT	Specificity	Uniform - Random	9.052	p<0.05
BRT	TSS	Uniform - BufferOut	0.2245	ns
BRT	TSS	Uniform - Random	8.2078	p<0.05
GAM	AUC	Uniform - BufferOut	2.4852	p<0.05
GAM	AUC	Uniform - Random	9.7106	p<0.05
GAM	CBI	Uniform - BufferOut	-2.9944	p<0.05
GAM	CBI	Uniform - Random	2.4044	p<0.05
GAM	RMSE	Uniform - BufferOut	-4.2491	p<0.05
GAM	RMSE	Uniform - Random	0.228	ns
GAM	Sensitivity	Uniform - BufferOut	2.7209	p<0.05
GAM	Sensitivity	Uniform - Random	5.3686	p<0.05
GAM	Specificity	Uniform - BufferOut	-0.5144	ns

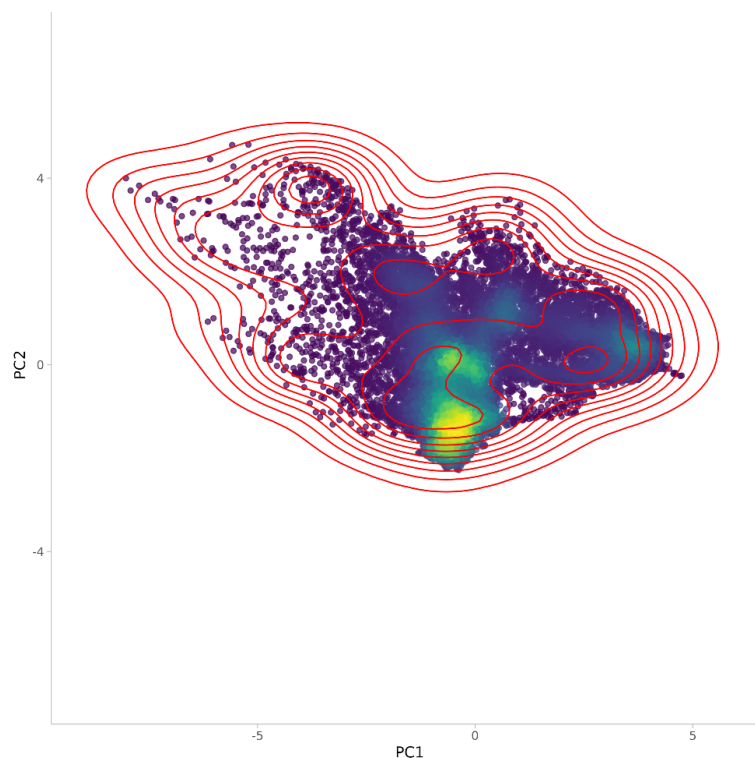
<b>Model</b>	<b>Metric</b>	<b>Comparisons</b>	<b>Z</b>	<b>P.val</b>
GAM	Specificity	Uniform - Random	8.4643	p<0.05
GAM	TSS	Uniform - BufferOut	0.0233	ns
GAM	TSS	Uniform - Random	8.2981	p<0.05
GLM	AUC	Uniform - BufferOut	-4.6005	p<0.05
GLM	AUC	Uniform - Random	-0.1257	ns
GLM	CBI	Uniform - BufferOut	0.726	ns
GLM	CBI	Uniform - Random	-5.4103	p<0.05
GLM	RMSE	Uniform - BufferOut	-3.1414	p<0.05
GLM	RMSE	Uniform - Random	2.7924	p<0.05
GLM	Sensitivity	Uniform - BufferOut	-1.5199	ns
GLM	Sensitivity	Uniform - Random	-2.8583	p<0.05
GLM	Specificity	Uniform - BufferOut	-5.522	p<0.05
GLM	Specificity	Uniform - Random	1.4241	ns
GLM	TSS	Uniform - BufferOut	-4.54	p<0.05
GLM	TSS	Uniform - Random	-0.3467	ns
Maxent	AUC	Uniform - BufferOut	2.4852	p<0.05
Maxent	AUC	Uniform - Random	9.7106	p<0.05
Maxent	CBI	Uniform - BufferOut	-7.9909	p<0.05
Maxent	CBI	Uniform - Random	0.4514	ns
Maxent	RMSE	Uniform - BufferOut	-2.8994	p<0.05
Maxent	RMSE	Uniform - Random	2.7528	p<0.05
Maxent	Sensitivity	Uniform - BufferOut	4.284	p<0.05
Maxent	Sensitivity	Uniform - Random	4.2468	p<0.05
Maxent	Specificity	Uniform - BufferOut	-0.8195	ns
Maxent	Specificity	Uniform - Random	7.6853	p<0.05
Maxent	TSS	Uniform - BufferOut	-0.1257	ns
Maxent	TSS	Uniform - Random	8.1398	p<0.05
RF	AUC	Uniform - BufferOut	2.6549	p<0.05



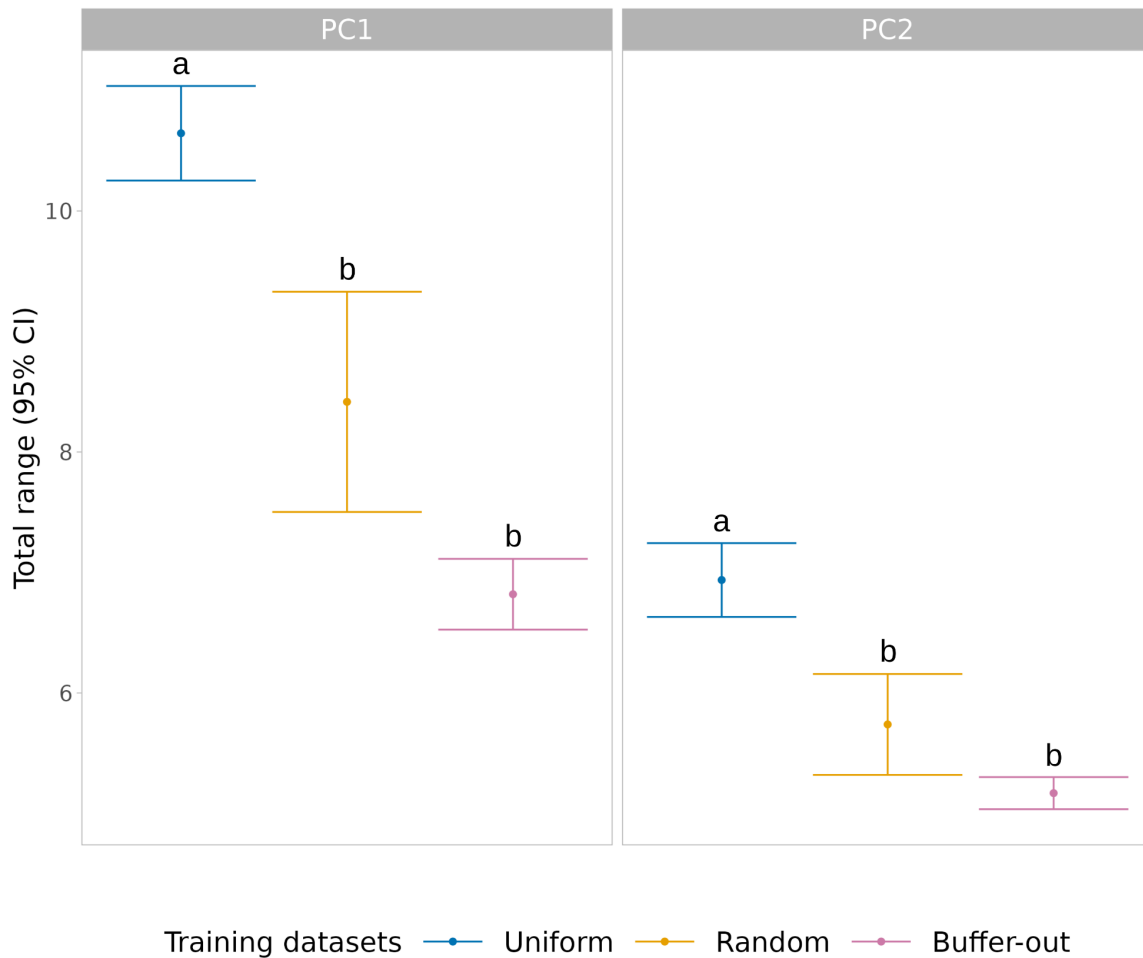
Model	Metric	Comparisons	Z	P.val
RF	AUC	Uniform - Random	9.7143	p<0.05
RF	CBI	Uniform - BufferOut	0.1289	ns
RF	CBI	Uniform - Random	8.4596	p<0.05
RF	RMSE	Uniform - BufferOut	-2.4619	p<0.05
RF	RMSE	Uniform - Random	3.1183	p<0.05
RF	Sensitivity	Uniform - BufferOut	-0.3741	ns
RF	Sensitivity	Uniform - Random	0.046	ns
RF	Specificity	Uniform - BufferOut	0.8738	ns
RF	Specificity	Uniform - Random	8.9689	p<0.05
RF	TSS	Uniform - BufferOut	0.0921	ns
RF	TSS	Uniform - Random	8.1664	p<0.05



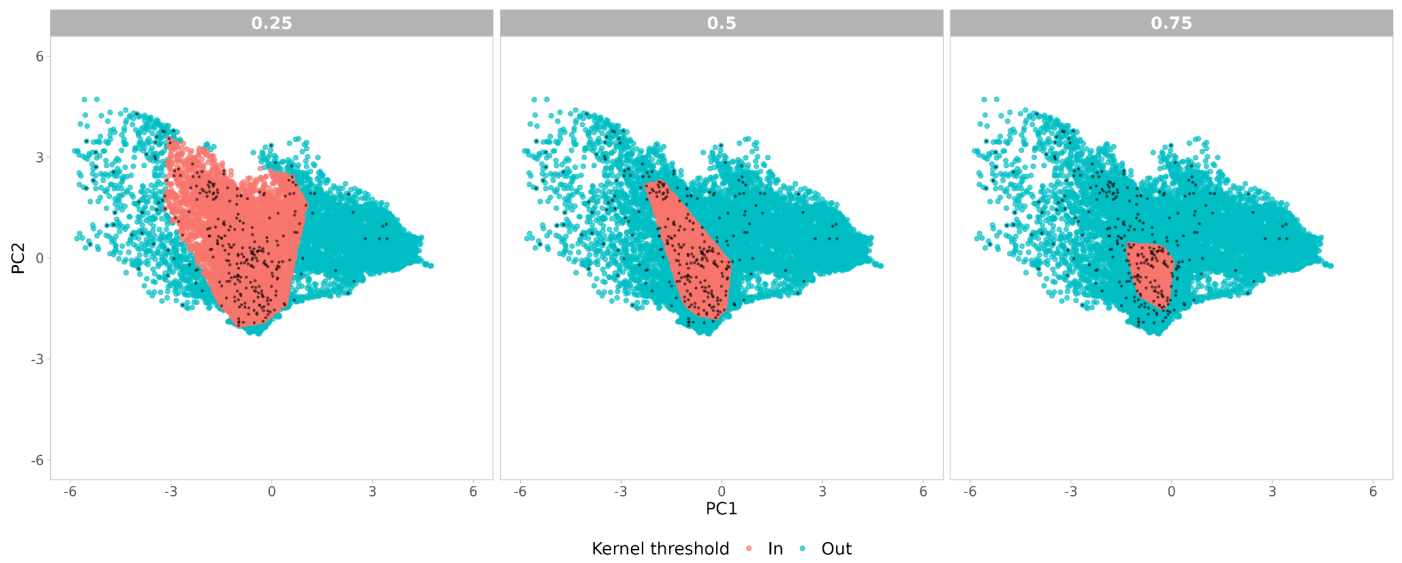
**S1.1:** Post-hoc multiple comparisons with one-tailed Dunn's rank sum test ( $\alpha = 0.05$ ; omnibus test was always significant with  $P < 0.05$ , data not shown). All the comparisons were performed assuming that the performance of the uniform sampling strategy was higher than the other two sampling strategies: a) relative proportion of the significant comparisons aggregated by sampling strategy; b) relative proportion of the significant comparisons aggregated by sampling strategy and metric.



**Figure S1.2:** Bivariate density plot of principal component scores associated with the pseudo-absences sampled for a virtual species using the uniform approach.



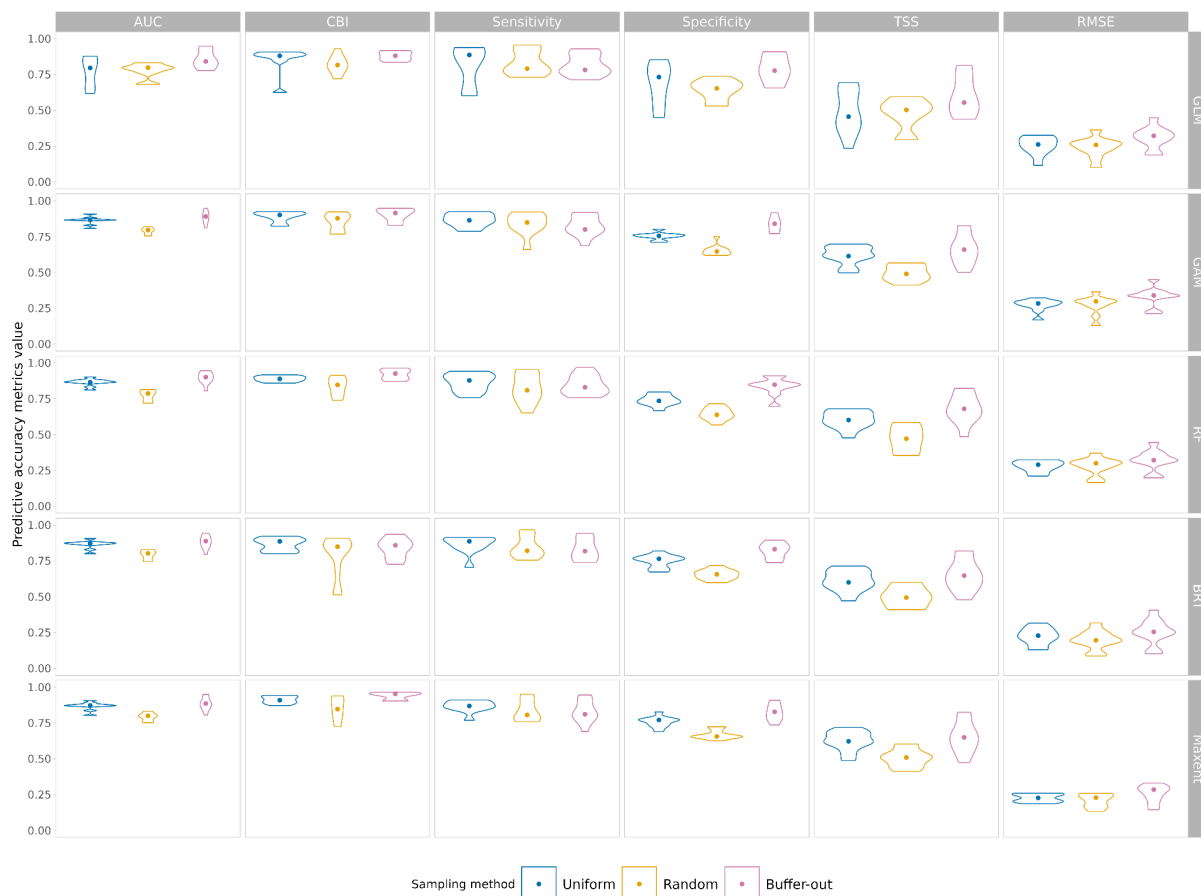
**Figure S1.3:** Mean (points) and 95% confidence interval (error bars) of the principal components total range (max PC-score - min PC-score) captured by the three sampling strategies. Two-tailed Kruskal-Wallis test for PC1:  $\chi^2 = 21.54$ ,  $df = 2$ ,  $p$ -value  $< 0.001$ ; Kruskal-Wallis test for PC2:  $\chi^2 = 14.91$ ,  $df = 2$ ,  $p$ -value  $< 0.001$ . Letters denote significant differences using Dunn's test, p-values were adjusted using Holm's correction. Colours are associated with the three sampling strategies used to sample the pseudo-absences (uniform in blue, random in yellow and buffer-out in pink).



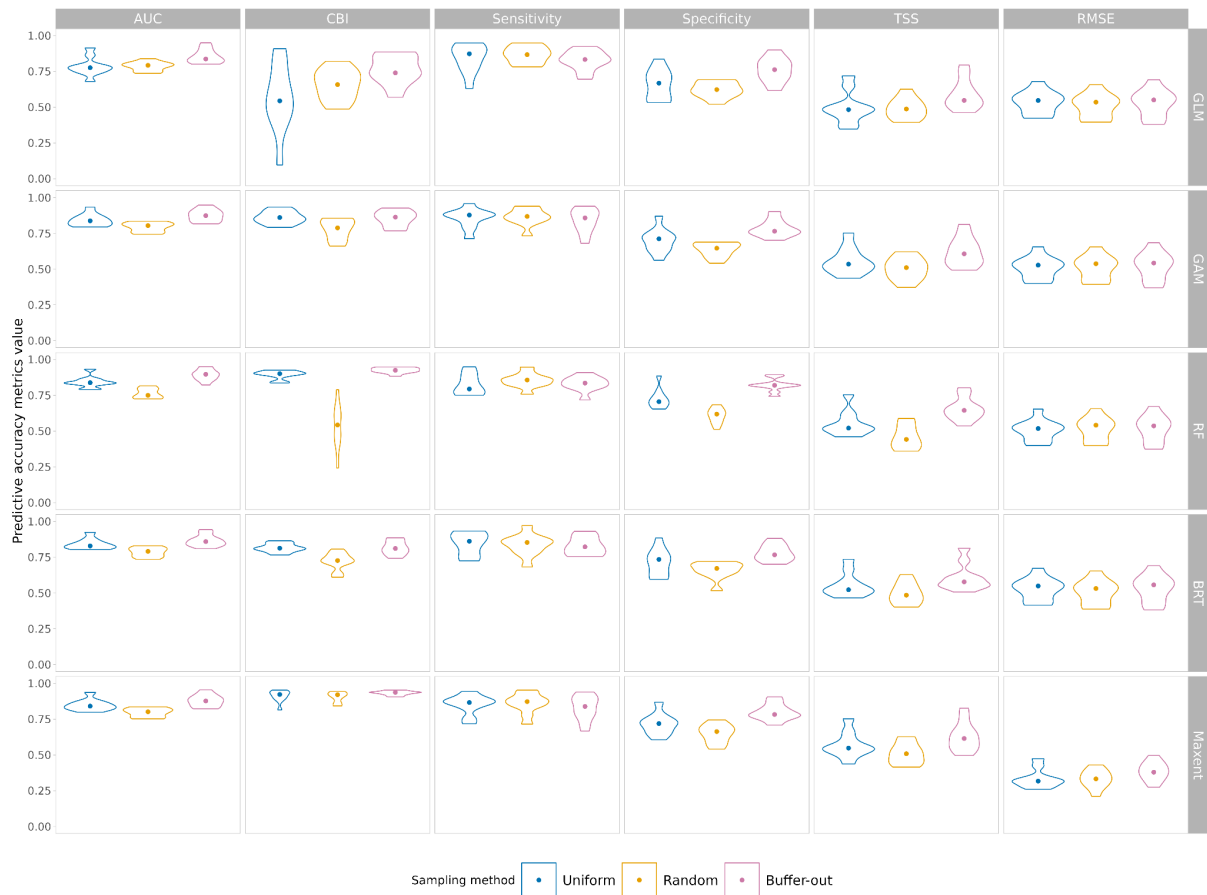
**Figure S1.4:** Effect of setting different kernel thresholds on the inclusion/exclusion of pseudo-absences eventually sampled using the uniform approach (black dots are the true virtual species presences represented within the environmental space). Setting a low value of the kernel threshold (e.g., 0.25) increases the portion of the environmental space excluded from the uniform sampling; in contrast, setting a high value of the kernel threshold increases the portion of the environmental space available for the uniform sampling.

## Supplementary Material 2

To test the potential effect of different sample prevalence values, we repeated the entire workflow on 10 virtual species using two different prevalence values: 0.5 and 0.1. In both cases, we obtained a dataset consisting of 300 presences, which we then combined with a second dataset of 600 (for sample prevalence 0.5) and 3,000 (for sample prevalence 0.1) pseudo-absences.

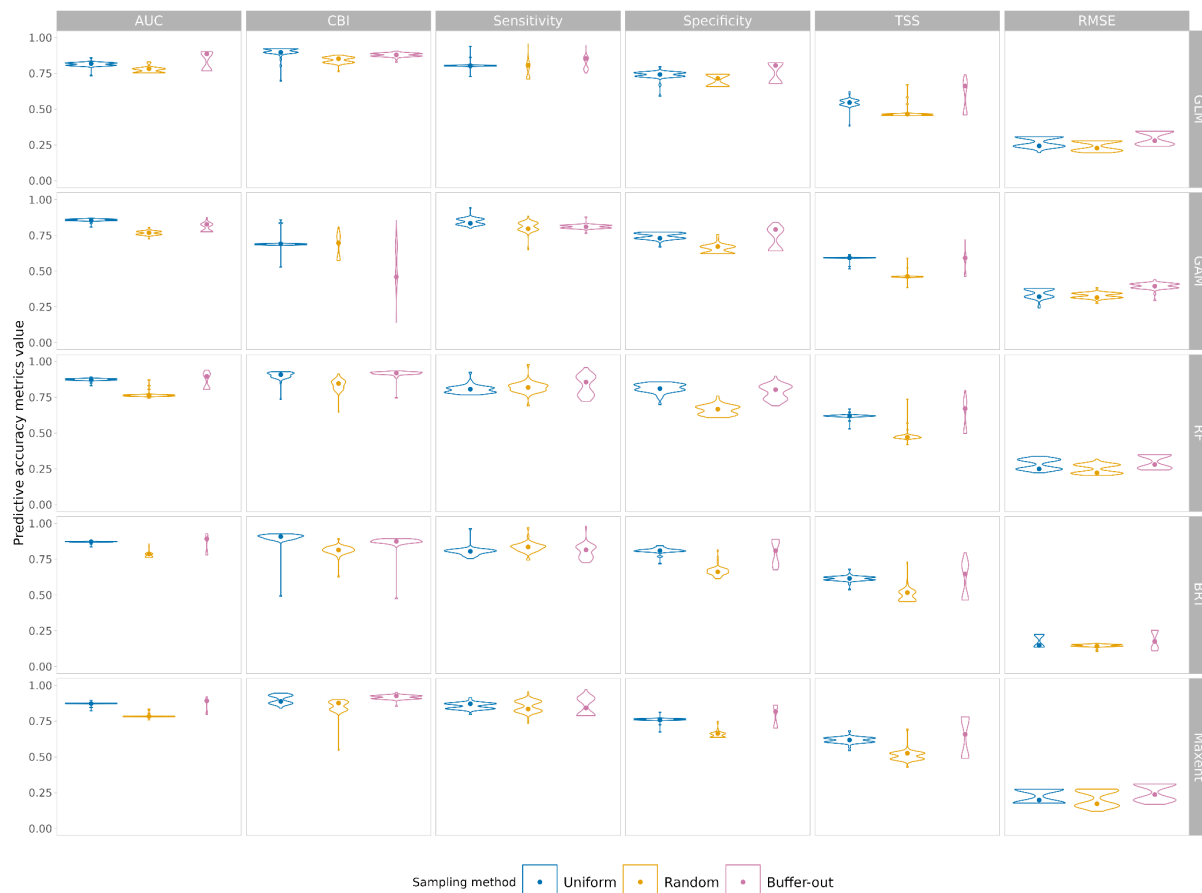


**Figure S2.1:** Violin plots reporting the distribution of the values of the metrics of predictive performance for the habitat suitability models of 10 virtual species (dots represent median values of the metrics of predictive performance), considering 5 predictors, and using a sample prevalence equal to 0.5. Columns indicate the different performance metrics, while rows are associated with the modelling algorithms used to fit the habitat suitability models.



**Figure S2.2:** Violin plots reporting the distribution of the values of the metrics of predictive performance for the habitat suitability models of 10 virtual species (dots represents median values of the metrics of predictive performance), considering 5 predictors, and using a sample prevalence equal to 0.1. Columns indicate the different performance metrics, while rows are associated with the modelling algorithms used to fit the habitat suitability models.

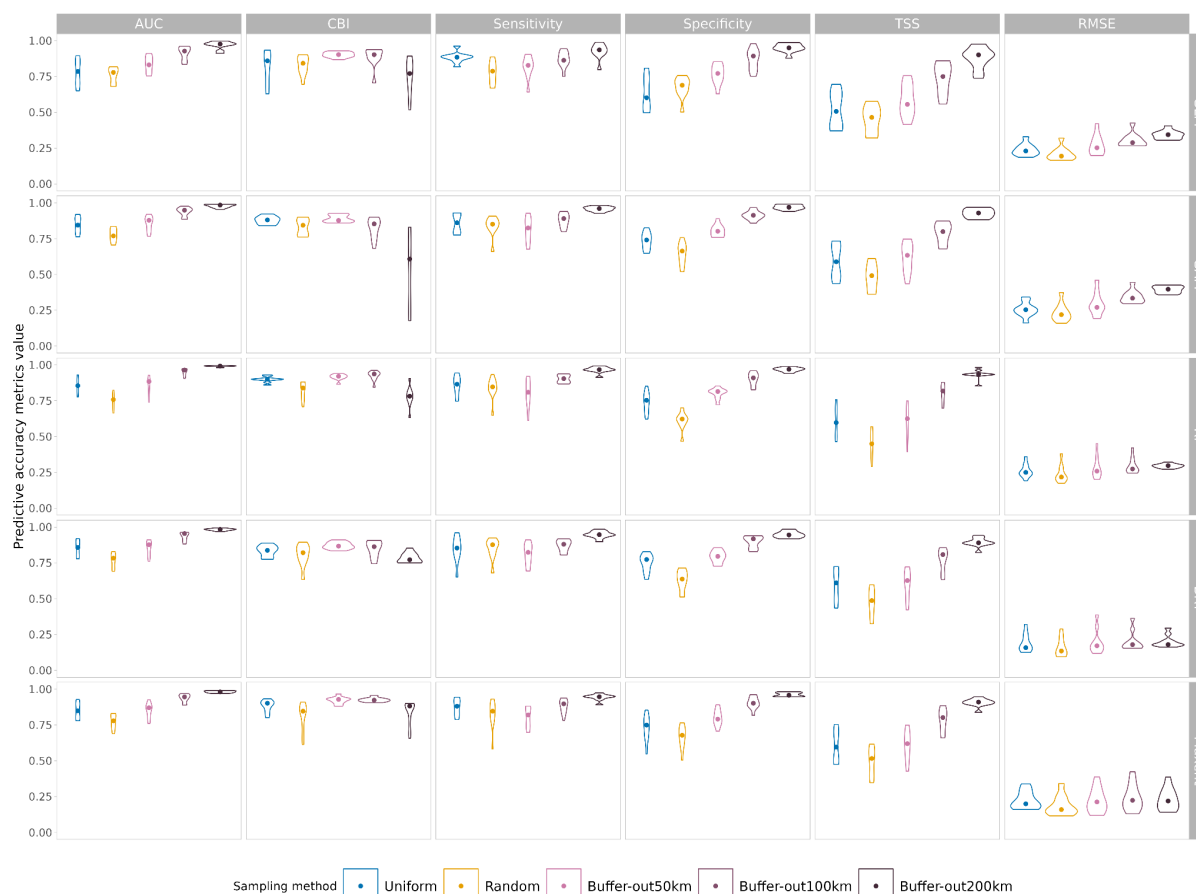
## Supplementary Material 3



**Fig. S3:** Violin plots reporting the distribution of the values of the metrics of predictive performance for the habitat suitability models of 50 virtual species modelled as a function of 19 bioclimatic predictors, and setting sample prevalence equal to 1 (i.e., same number of presences and pseudo-absences). Dots represent median values of the metrics of predictive accuracy. Columns indicate the different performance metrics, while rows are associated with the modelling algorithms used to compute HSMs. AUC = Area Under the Curve; CBI = Continuous Boyce Index, TSS = True Skill Statistic; RMSE = Root Mean Squared Error; GLM = generalised linear model; GAM = generalised additive model; RF = random forest; BRT = boosted regression trees.

## Supplementary Material 4

To test the potential effect of different sizes of the buffer sizes on the buffer-out approach, we repeated the entire workflow on 10 virtual species with three different radius lengths: 50, 100 and 200 km. We kept the training dataset with a sample prevalence equal to 1, consisting of 300 presences and 300 pseudo-absences.



**Fig. S4.1:** Violin plots reporting the distribution of the values of the metrics of predictive performance for the habitat suitability models of 10 virtual species modelled as a function of 5 bioclimatic predictors, and setting sample prevalence equal to 1 (i.e., same number of presences and pseudo-absences). We varied the size of the radius for the buffer-out approach, setting it to/using 50, 100 and 200 km. Dots represent median values of the metrics of predictive accuracy. Columns indicate the different performance metrics, while rows are associated with the modelling algorithms used to fit the habitat suitability models. AUC = Area Under the Curve; CBI = Continuous Boyce Index, TSS = True Skill Statistic; RMSE = Root Mean Squared Error; GLM = generalised linear model; GAM = generalised additive model; RF = random forest; BRT = boosted regression trees.



## Supplementary Material 5

### Case study on the realised distribution of *Fagus sylvatica* in Western Europe

#### Methods

To illustrate how to apply the uniform approach with the `USE` R package, we modelled the realised distribution of *Fagus sylvatica* in Italy, France and Spain (hereafter, western Europe). We chose *F. sylvatica* as an example species because its distribution and biogeographic history is well-known across Europe (Magri et al., 2006; Poli et al., 2022). For the sake of simplicity, we restricted the area of investigation to western Europe and used two modelling algorithms. Indeed, the case study of *F. sylvatica* is only intended as a practical example to show how the `USE` package operates, while not providing a further comparison on the predictive performance of HSMs fitted on data collected through different sampling strategies (as already done using virtual species, see main manuscript). We gathered data on the presence of *F. sylvatica* from the open EU-Forest dataset (Mauri et al., 2017), which compiles observations on European tree species from national inventories and other similar sources (see Mauri et al., 2017 for further information about EU-Forest). EU-Forest data consist of presence records of tree species exhaustively collected across Europe, and then aggregated to a  $1 \times 1$  km resolution grid. This lets us assume with a certain degree of confidence that the EU-Forest dataset provided a geographically unbiased sample of presence records for *F. sylvatica* in western Europe.

Across our study area, the EU-Forest dataset included a total of 12,444 presence records for *F. sylvatica*, which we sub-sampled within the environmental space to retrieve both a training and a testing (for internal validation) presence dataset. To this aim, we generated a 2-dimensional environmental space using all 19 bioclimatic variables available from WorldClim. Then, we used the function `USE::uniformSampling` to uniformly sample presence records within the environmental space. Note that this approach is conceptually similar to the spatial-thinning proposed by Aiello-Lammens et al. (2015), which aims at reducing the clustering of presences within the geographical space (Sillero and Barbosa, 2020), except that here we applied it within the environmental space. The obtained training and testing presence datasets were then combined to obtain the training and testing pseudo-absence datasets using the `paSampling`

function from the USE package. In particular, all presence records available for *Fagus sylvatica* were used to recover the core area of the species' bioclimatic niche within the environmental space. This allowed filtering out the pseudo-absences likely associated with suitable locations for the species (see step 1 in section 2.2.1 in the main text). The final sample size of the pseudo-absences included in the training and testing (internal validation) datasets were 1,826 and 991, respectively. Note that the sample size of the presence data included in the training and testing datasets were 1,827 and 991, respectively. Also note that prevalence was fixed to approx. 1 in both the training and testing dataset.

Finally, we derived a completely independent testing (external validation) dataset using presence and true absence data from sPlotOpen (Sabatini et al., 2021). The sPlotOpen dataset is an open-access subset of sPlot, one of the most comprehensive global databases of vegetation records (Sabatini et al., 2021). Here, we used sPlotOpen to gather *F. sylvatica* presences ( $n = 367$ ), and to derive true absence data from those vegetation plots where *F. sylvatica* was not recorded ( $n = 4,162$ ). As done for the EU-Forest dataset, we considered only sPlotOpen vegetation plots occurring in western Europe (i.e., Italy, France and Spain).

Then, we modelled the realised distribution of *F. sylvatica* as a function of a set of WorldClim bioclimatic variables. For simplicity, we solely focused on the climatic niche of *Fagus sylvatica*, although we acknowledge that other drivers than climate equally contribute in shaping the distribution of this species, especially so at local scales (Mellert et al., 2018). As modelling techniques, we used a 'logit' link binomial generalised linear model (binomial GLM) and random forests (RF, fitted using `ranger::ranger`; Wright and Ziegler, 2017). To reduce multicollinearity, we selected a subset of the 19 bioclimatic variables using the R function `caret::findCorrelation` function (Kuhn, 2021) (setting the pairwise-correlation threshold to 0.6). The bioclimatic variables eventually kept to fit the HSM for *F. sylvatica* were: BIO6 (minimum temperature of the coldest month); BIO7 (temperature annual range); BIO8 (mean temperature of the wettest quarter). Also, we used the latitudinal position of the presence and pseudo-absence records (hereafter, latitude) as an additional predictor to account for the effect of factors affecting the latitudinal gradient of the distribution of *F. sylvatica* that were not included in the model. An example of such factors is the

species biogeographic history of post-glacial recolonization towards northern Europe (Magri et al., 2006). To account for non-linearity in the profile of Pearson's residuals and improve the fit of the binomial GLM, we introduced second order polynomial terms for BIO6, BIO7 and latitude. The predictive performance of the fitted models was assessed on three different types of data: (i) the (internal) testing dataset derived from the EU-Forest dataset; (ii) 5 partitions of the training dataset (i.e., a 5-fold cross-validation); and (iii) the independent (external) testing dataset derived from sPlotOpen. As predictive accuracy metrics, we used the true skill statistics (TSS) and the continuous Boyce index (CBI). A TSS value greater than 0.5 is often considered to indicate good predictions. Positive values of CBI indicate that presences predicted by the model are consistent with the distribution of presences in the testing dataset. On the contrary, TSS and CBI values close to zero indicate that the model does not perform differently from a model that randomly predicts presences and absences. Finally, negative values of the CBI indicate counter predictions, i.e., predicting low suitability in areas with high density of presence records (Hirzel et al. 2006).

Beyond model predictive metrics, we computed the following measures of goodness-of-fit: Tjur's  $R^2$  for the binomial GLM and the  $R^2$  for the RF.

A full description of the modelling procedure (from the sub-sampling of the presences and the collection of pseudo-absences to the assessment of the model predictive performance) can be found at [https://github.com/danddr/USE\\_paper/tree/main/Example](https://github.com/danddr/USE_paper/tree/main/Example).

## Results

Both the binomial GLM and the RF for *F. sylvatica* showed high predictive performances, regardless of the dataset used for testing (Table S5.1). Concerning the binomial GLM, the TSS was always equal to or above 0.41, with the lowest value obtained for the sPlotOpen testing dataset (0.41) and the highest for the EU-Forest dataset (0.61). Similarly, the lowest CBI was scored for the sPlotOpen dataset (0.88), while the highest for the EU-Forest dataset (0.99).

We obtained comparable results for the RF, with the lowest TSS obtained when using sPlotOpen as a testing dataset (0.52), while the EU-Forest dataset and the (average across) 5-fold cross validation resulted in TSS equal to 0.79 and 0.77, respectively. With respect to the CBI, the highest value was observed

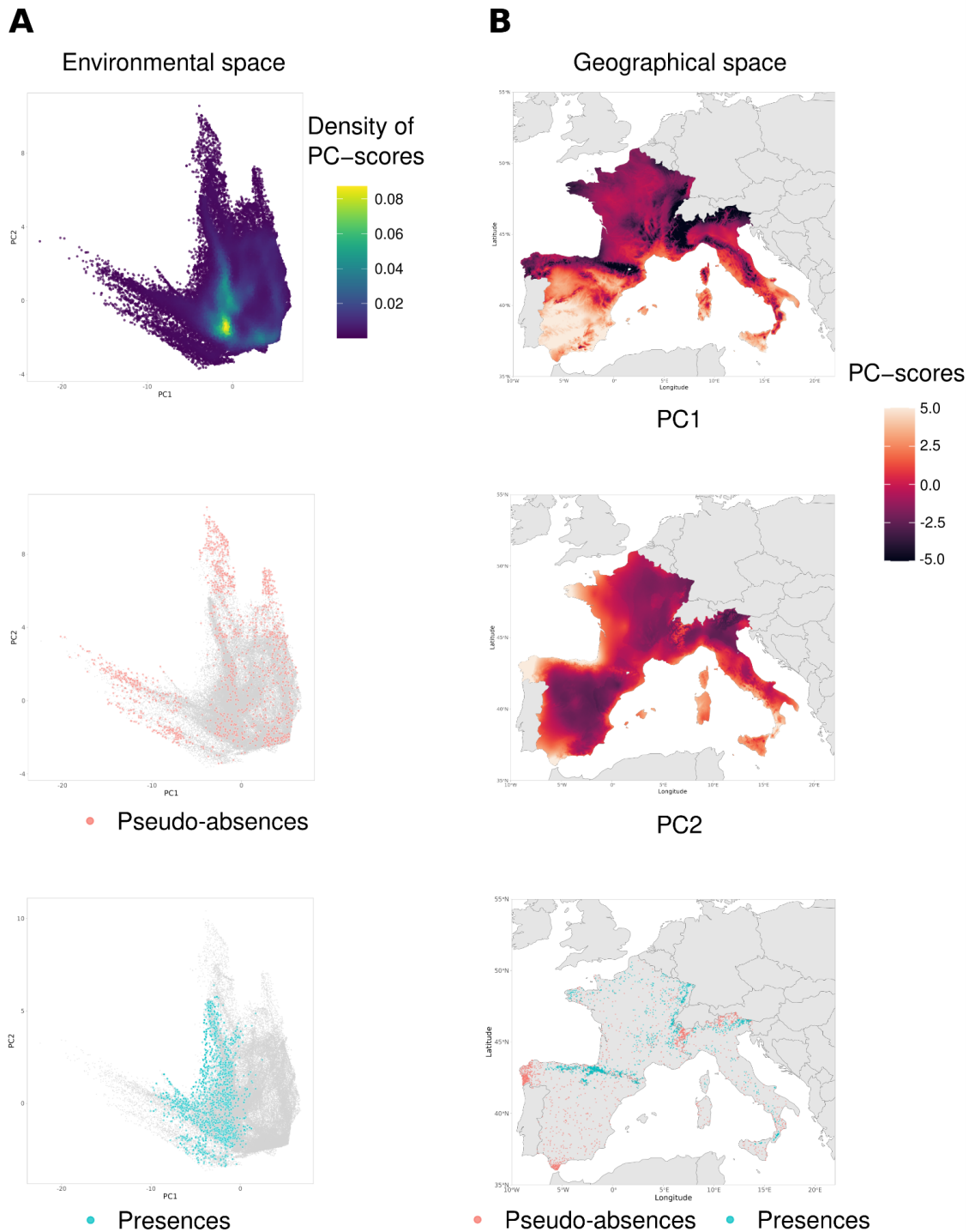
for the EU-Forest dataset (0.99), while the lowest was obtained using the sPlotOpen dataset (0.93).

Goodness-of-fit measures seemed to be affected by the modelling technique, with the  $R^2$  of the RF being 0.66, and the Tjur's  $R^2$  for the GLM being 0.36 (Tab. S5.1).

The pseudo-absences of *F. sylvatica* collected using the uniform approach were homogeneously distributed within the environmental space (Fig. S5.2a).

**Table S5.1:** Results of the habitat suitability models for *Fagus sylvatica* (generalised linear model, GLM, and random forest, RF). Models' predictive performance was assessed through internal (5-fold cross-validation and EU-Forest) and external (sPlotOpen) validation. TSS: true skill statistics; CBI: continuous Boyce index; R-sq: Tjur's  $R^2$  for the GLM, and  $R^2$  for RF. Values of TSS and CBI for the 5-fold cross-validation represent averages.

Validation dataset	GLM			RF		
	TSS	CBI	Tjur's $R^2$	TSS	CBI	$R^2$
5-fold CV	0.52	0.93		0.77	0.97	
EU-Forest	0.61	0.99	0.36	0.79	0.99	0.66
sPlotOpen	0.41	0.88		0.52	0.93	



**Figure S5.2:** (A) environmental space available for *Fagus sylvatica* in Italy, Spain and France, and the position of presences (light blue) and pseudo-absences (red) sampled within the environmental space using the uniform approach; (B) distribution of principal component scores across the geographical space, and location (across western Europe) of presences (light blue) and pseudo-absences (red) sampled using the uniform approach.

## References

- Aiello-Lammens, M. E., Boria, R. A., Radosavljevic, A., Vilela, B., and Anderson, R. P. (2015). sptin: an R package for spatial thinning of species occurrence records for use in ecological niche models. *Ecography*, 38(5):541–545.
- Hirzel, A. H., Le Lay, G., Helfer, V., Randin, C., and Guisan, A. (2006). Evaluating the ability of habitat suitability models to predict species presences. *Ecological modelling*, 199(2), 142-152.
- Kuhn, M. (2021). *caret: Classification and Regression Training*. R package version 6.0-88.
- Magri, D., Vendramin, G. G., Comps, B., Dupanloup, I., Geburek, T., Gömöry, D., ... and De Beaulieu, J. L. (2006). A new scenario for the Quaternary history of European beech populations: palaeobotanical evidence and genetic consequences. *New phytologist*, 171(1), 199-221.
- Mauri, A., Strona, G., and San-Miguel-Ayanz, J. (2017). Eu-forest, a high-resolution tree occurrence dataset for europe. *Scientific data*, 4(1):1–8.
- Mellert et al. (2018) Soil water storage appears to compensate for climatic aridity at the xeric margin of European tree species distribution. *European Journal of Forest Research*, 137: 79-92.
- Poli et al. (2022) Coupling fossil records and traditional discrimination metrics to test how genetic information improves species distribution models of the European beech *Fagus sylvatica*. *European Journal of Forest Research*, 141: 253–265
- Sabatini, F. M., Lenoir, J., Hattab, T., Arnst, E. A., Chytrý, M., Dengler, J., De Ruffray, P., Hennekens, S. M., Jandt, U., Jansen, F., et al. (2021). splotopen—an environmentally balanced, open-access, global dataset of vegetation plots. *Global Ecology and Biogeography*.
- Sillero, N. and Barbosa, A. M. (2020). Common mistakes in ecological niche models. *International Journal of Geographical Information Science*, pages 1–14.
- Wright, M. N. and Ziegler, A. (2017). ranger: A fast implementation of random forests for high dimensional data in C++ and R. *Journal of Statistical Software*, 77(1):1–17.