

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^{1,*†}, Enrico Tordoni^{2†}, Jonathan Lenoir³,
5 Jonas J. Lembrechts⁴, Sophie O. Vanwambeke¹,
6 Duccio Rocchini^{5,6}, and Manuele Bazzichetto^{7†}

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

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

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

13 ⁴ Research Group Plants and Ecosystems, University of Antwerp, Belgium

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

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

18 ⁷ 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

27

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 set.
- 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 geographic 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 realized 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
93 scientific literature (Sillero and Barbosa, 2020), but they 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 use the term pseudo-
101 absences for simplicity and alignment with our study.

102 The most common approaches for sampling pseudo-absences involve (i) randomly
103 surveying a large number of points across the study area (e.g., 10,000; Barbet-Massin et
104 al., 2012; Iturbide et al., 2015; Støa et al., 2019, Hysen et al., 2022) or (ii) sampling them

USE: a novel approach to uniformly sample the environmental space

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

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

USE: a novel approach to uniformly sample the environmental space

130 of HSMs and it is particularly critical for machine learning techniques, while regression
131 techniques such as Generalised Linear Models seem to be less affected (Barbet-Massin et
132 al., 2012; Grime, Whittaker and Hortal, 2020; Valavi et al., 2021). So far, class overlap has
133 been addressed using resampling techniques more oriented to adjusting an unbalanced
134 number of classes in the response variable (i.e., the ‘up-’ or ‘down-sampling’ approach;
135 Valavi et al., 2021), irrespective of the technique to obtain pseudo-absences.

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

154

USE: a novel approach to uniformly sample the environmental space

155 2 Methods

156 2.1 Simulation of virtual species

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

173 2.2 Sampling of the pseudo-absences

174 Regardless of the sampling approach and modelling technique used to calibrate the HSMS,
175 the ratio between the number of presences and pseudo-absences in the training datasets

USE: a novel approach to uniformly sample the environmental space

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

182 *2.2.1 Uniform approach: pseudo-absences sampled within the environmental space*

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

- 197 1. First, kernel density estimation (a statistical technique used to estimate the underlying
198 probability distribution of a set of data points by smoothing them with a kernel

USE: a novel approach to uniformly sample the environmental space

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

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

- 220 - Within each cell of the sampling grid, the average (squared) Euclidean distance
221 between the pseudo-absences (PC-scores) in the cell and the centroid of their
222 convex hull is computed;

USE: a novel approach to uniformly sample the environmental space

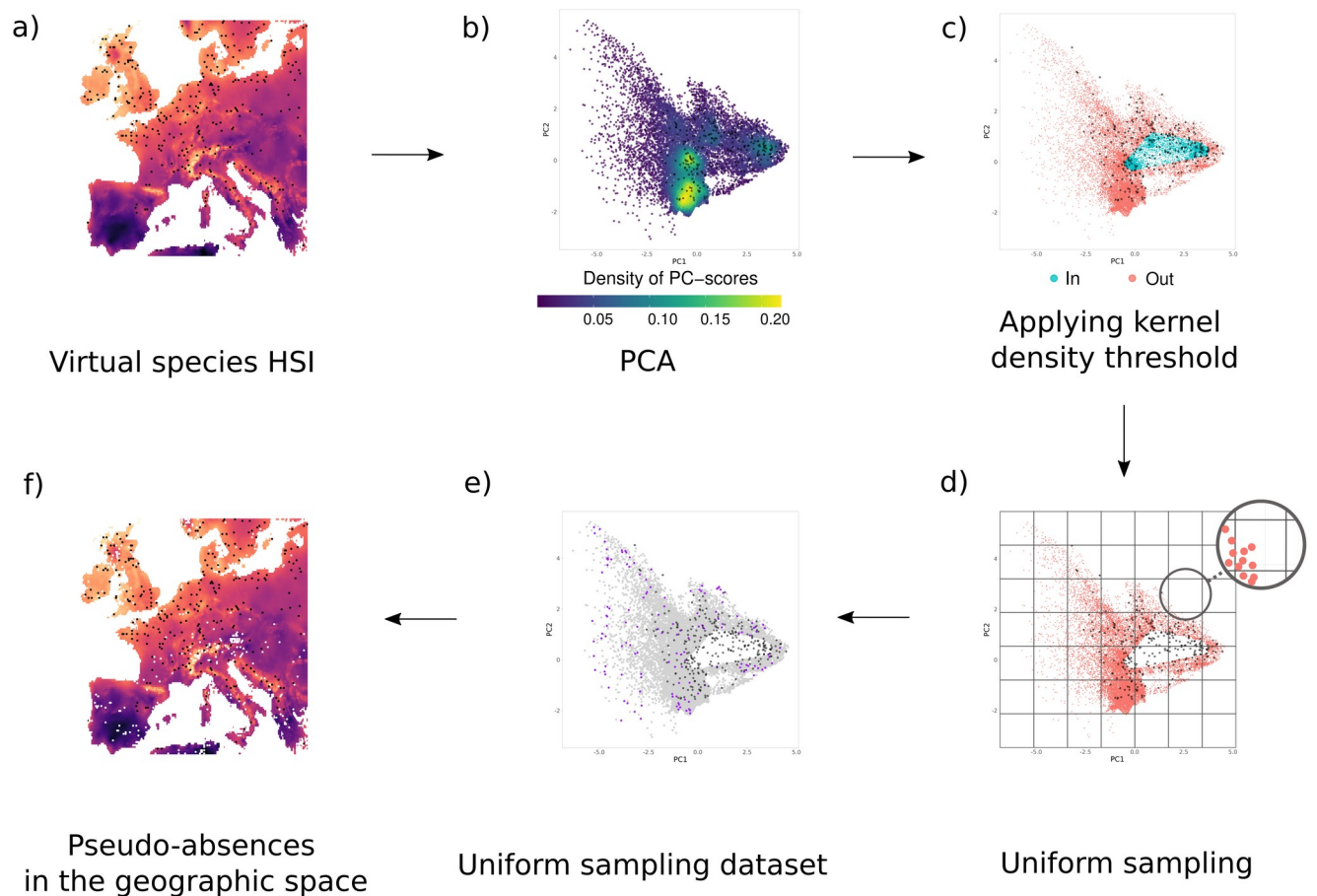
- 223 - Once this metric is computed across all cells of the sampling grid, the average mean
224 value is computed across all cells (hereafter, grid average);
- 225 - The procedure above is separately repeated on different sampling grids of
226 increasing resolution (i.e., increasing number of cells);
- 227 - The resulting set of grid averages (one per resolution) are used as a measure of the
228 aggregation among pseudo-absences within the cells of the sampling grids. This
229 value is compared across resolutions and the best grid is chosen as the one
230 providing the best trade-off between resolution and average distance among points
231 within cells (i.e., resolution that allows uniformly sampling the environmental space
232 without overfitting it). More specifically, the best grid is the one whose resolution is
233 just below that which would not allow the average distance among pseudo-absences
234 to be reduced by more than 10% (other values can be set by the user).

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

242 The total number of pseudo-absences sampled within each cell of the sampling grid can be
243 set by the user (using the argument `n.tr`, default `n.tr = 5`), who can also indicate a
244 desired sample prevalence. If the sample prevalence is not specified, fewer pseudo-absences
245 are likely to be eventually sampled than expected (i.e., `n.tr × number of cells`). This happens
246 because (i) no pseudo-absence points are collected in empty cells, and (ii) less pseudo-

USE: a novel approach to uniformly sample the environmental space

247 absence points than `n.tr` are available within the cells at the boundary of the environmental
248 space (see zooming window in Figure 1d). Similarly, no pseudo-absences are collected within
249 the core area of the presences (excluded in step 1). If a sample prevalence is set by the user,
250 the sampling grid is surveyed until the chosen sample prevalence is reached by the
251 algorithm.



252 **Figure 1:** Flowchart representing the step-by-step procedure for implementing the uniform
253 approach: a) habitat suitability index (HSI) of the i -th virtual species (VS; lighter colours
254 indicate higher habitat suitability and black dots represent presence points in the
255 geographical space); b) PCA performed on the environmental variables in the study region

USE: a novel approach to uniformly sample the environmental space

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

270 *2.2.2 Pseudo-absences sampled within the geographical extent*

271 The sampling of pseudo-absences within the geographical extent was conducted using the
272 random and buffer-out approaches. For the random approach (Barbet-Massin et al. 2012;
273 Iturbide et al., 2015; Støa et al., 2019), we simply generated 300 random pseudo-absences
274 across the studied geographical extent. For the buffer-out approach (Bedia et al., 2013), we
275 created a buffer of 50 km radius around each presence location, and we then randomly
276 sampled pseudo-absences outside the presence-specific buffers, but within the convex hull
277 of the species geographical distribution (i.e., the convex hull that connects the outer
278 presences of the species and thus delimits the range actually covered by the species in the
279 geographical space). To test the effect of the length of the radius on the buffer-out

USE: a novel approach to uniformly sample the environmental space

280 approach, we performed a sensitivity analysis on 10 VS using the following radius lengths:
281 50, 100 and 200 km.

282 2.3 Comparison among sampling strategies

283 2.3.1 Predictive performance comparison

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) in 5 training (70% observations) and testing (30%) sets,
290 which we used to calibrate and validate five modelling algorithms: (i) binomial generalised
291 linear models with 'logit' link (GLMs); (ii) generalised additive models (GAMs); (iii) random
292 forests (RFs); (iv) boosted regression trees (BRTs); and (v) MaxEnt. In total, we fitted 3,750
293 HSMs (50 VS species \times 3 different sets of pseudo-absences \times 5 modelling algorithms \times 5
294 replicates of 70-30% partitions). To fit the HSMs, we used the R package `sdm` (Naimi and
295 Araújo, 2016). Although we acknowledge the importance of fine-tuning HSMs (Fourcade,
296 2021), we kept model settings at their default value since it would have been unfeasible to
297 individually parametrise each algorithm for all 50 VS and sampling strategies. A detailed
298 representation of the workflow of the analyses is shown in Fig. 2. Furthermore, we
299 acknowledge that our use of MaxEnt did not conform with the general recommendations for
300 its adequate implementation (e.g., using 10,000 background points; Cobos et al., 2019;
301 Kaas et al., 2021). Nonetheless, we included it in the comparison of models' performance
302 due to its wide usage within the HSMs community.

303 After fitting HSMs for all the 50 VS, we compared the predictive performance

USE: a novel approach to uniformly sample the environmental space

304 associated with each combination of sampling approaches and modelling techniques by
305 computing the following metrics: (i) Area Under the receiver operating characteristic Curve
306 (AUC); (ii) Continuous Boyce Index (CBI); (iii) sensitivity; (iv) specificity; (v) True Skill
307 Statistics (TSS); and (vi) Root Mean Squared Error (RMSE). The RMSE was computed by
308 comparing the true (i.e., simulated) habitat suitability of the focal VS against the one
309 predicted by each combination of modelling and sampling approach. A detailed description
310 of the above-mentioned modelling techniques and validation metrics can be found in
311 Guisan et al. (2017). To compare the predictive performance of the HSMs fitted under
312 different combinations of sampling strategy and modelling technique, we visually assessed
313 the results of the 50 VS simulations using violin plots reporting the distribution of the values
314 of the predictive performance metrics listed above. Furthermore, we tested for statistical
315 differences among the predictive performance of the sampling strategies using Kruskal-
316 Wallis tests, followed by Dunn's post hoc rank sum comparisons using the `dunn.test` R
317 package (Dinno, 2017) (p-values for multiple comparisons adjusted using Holm correction).

318 To test the potential effect on our comparison of varying the number of bioclimatic
319 variables used, we repeated the entire workflow on 50 VS using all the 19 bioclimatic
320 variables for both the selection of pseudo-absences and as predictors in the five HSM
321 algorithm. To test the potential effect on our comparison of varying sample prevalence, we
322 repeated the entire workflow on 10 VS using two additional prevalence values, namely 0.5
323 and 0.1. Specifically, for each VS, we generated two additional training datasets with 300
324 presences, but we combined them with 600 and 3,000 pseudo-absences, to achieve
325 sample prevalence of 0.5 and 0.1 respectively.

326 *2.3.2 Sample location bias and class overlap*

327 To assess the intensity of sample location bias associated with the different sampling

USE: a novel approach to uniformly sample the environmental space

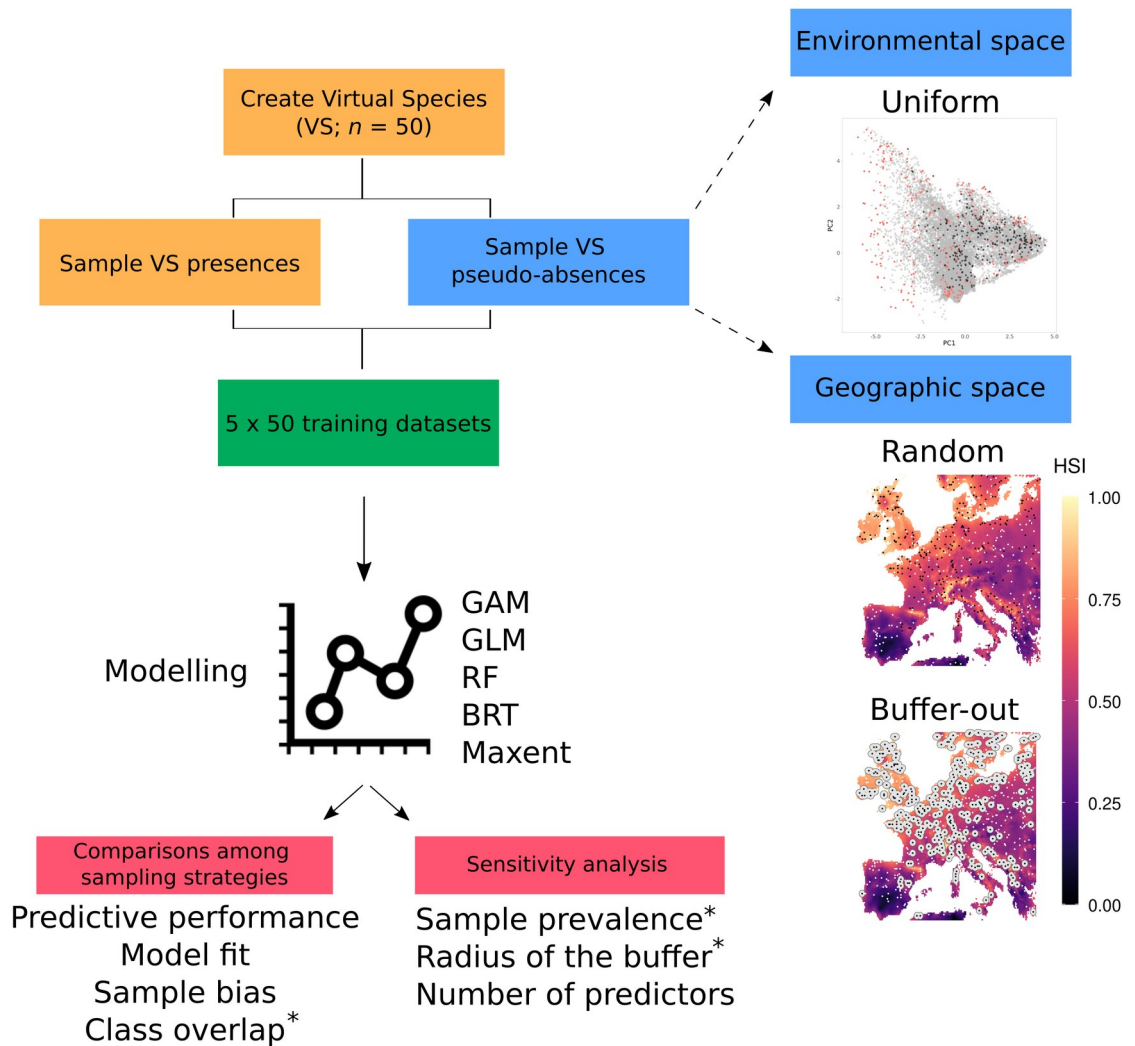
328 strategies, we extracted the pseudo-absences of a single VS and map their aggregation
329 within the environmental space using bivariate density plots. The aim was to identify which,
330 among the three sampling strategies, was more subject to oversampling particular
331 environmental conditions within the geographical space. In principle, the sampling
332 strategies more affected by sample location bias would exhibit a clear aggregation of
333 pseudo-absences within the environmental space. We visually assessed the areas of the
334 environmental space sampled by the different sampling strategies using the function
335 `geom_density_2d` of the `ggplot2` R package (Wickham, 2016). This function performs a
336 2D kernel density estimation using the `kde2d` function of the `MASS` R package (Venables
337 and Ripley, 2002) and displays the results with contours.

338 To assess the effectiveness of the uniform approach for mitigating class overlap, we
339 simulated 10 new VS, sampled their presences and pseudo-absences using the three
340 sampling strategies and mapped the position of the presence and pseudo-absence points
341 within the environmental space following the procedure explained in section 2.2.1 and
342 Figure 1a,b. Then, we computed the Gaussian hypervolume of the presences and pseudo-
343 absences using the `hypervolumes` R package (Blonder, 2022), and calculated the overlap
344 between them. Statistically significant differences in the degree of overlap were tested using
345 one-way ANOVA and Tukey HSD test.

346 2.4 Real-case study

347 To illustrate how to apply the uniform approach with the `USE` R package, we modelled the
348 realised distribution of *Fagus sylvatica* in Italy, France and Spain. We chose *F. sylvatica* as
349 a target species because its distribution and biogeographic history is well-known across
350 Europe (Magri et al., 2006; Poli et al., 2022). The whole analysis of *F. sylvatica* is described
351 in S5, and the R code to replicate it can be found at: https://github.com/danddr/USE_paper.

352



353 **Figure 2** Overall workflow of the analysis described in the Methods section. The '*' is
 354 associated with analyses (i.e., class overlap, sample prevalence, radius of the buffer)
 355 performed on $n = 10$ VS.

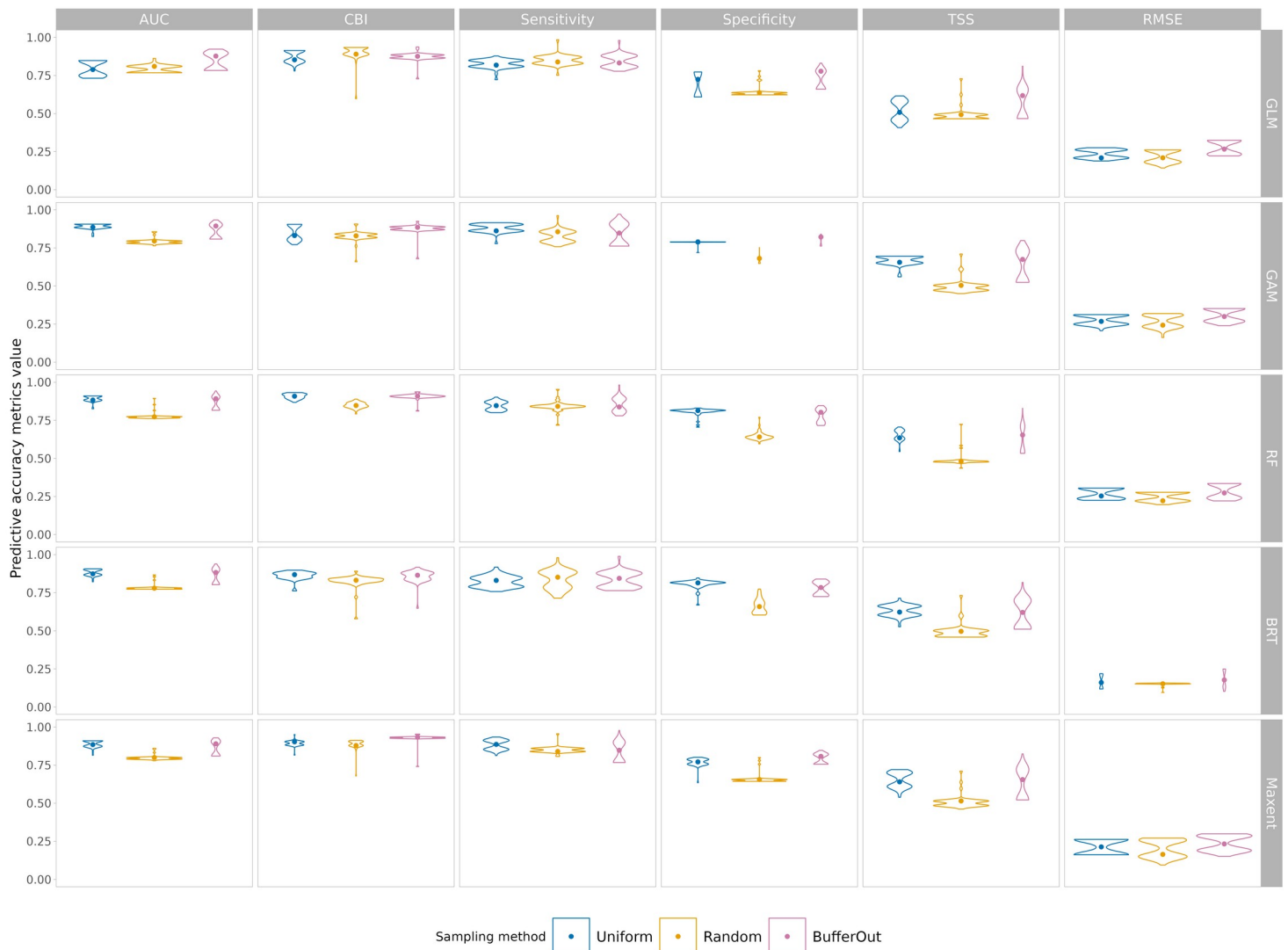
USE: a novel approach to uniformly sample the environmental space

356 3 Results

357 3.1 Comparison of the predictive performance associated with geographical vs 358 environmental sampling

359 Overall, the uniform approach performed equal to or better than the geographical
360 approaches in terms of out-of-sample prediction (Fig. 3). Pairwise comparisons between the
361 predictive accuracy performance of the random and buffer-out approaches against the
362 uniform approach showed statistically significant differences in 73% and 53% of the
363 combinations, respectively. However, these differences were algorithm- and metric-
364 dependent and did not point to a higher predictive performance of the uniform approach
365 (Fig. 3, Tab. S1, Fig. S1.1). The pattern of the differences among predictive performance
366 metrics was consistent among prevalence values (Fig. S2.1-2.2) and number of bioclimatic
367 variables used in the models (Fig. S3). Increasing the buffer radius length (Fig. S4),
368 resulted in higher predictive performance of the buffer-out approach for some metrics (AUC,
369 TSS, Specificity), while for CBI, Sensitivity and RMSE results remained comparable with
370 those presented in Fig. 3.

USE: a novel approach to uniformly sample the environmental space



371 **Figure 3:** Violin plots reporting the distribution of the values of the metrics of predictive
 372 performance for the HSMs of the 50 VS, as modelled using 5 randomly selected bioclimatic
 373 predictors and setting sample prevalence equal to 1 (i.e., same number of presences and pseudo-
 374 absences). Dots represent median values of the metrics of predictive accuracy. Columns indicate
 375 the different performance metrics, while rows are associated with the modelling techniques used to
 376 fit the HSMs. Higher values in all metrics but RMSE reflect higher predictive performance. AUC =
 377 Area Under the Curve; CBI = Continuous Boyce Index, TSS = True Skill Statistic; RMSE = Root
 378 Mean Squared Error; GLM = Generalized Linear Model; GAM = Generalized Additive Model; RF =
 379 Random Forest; BRT = Boosted Regression Trees.

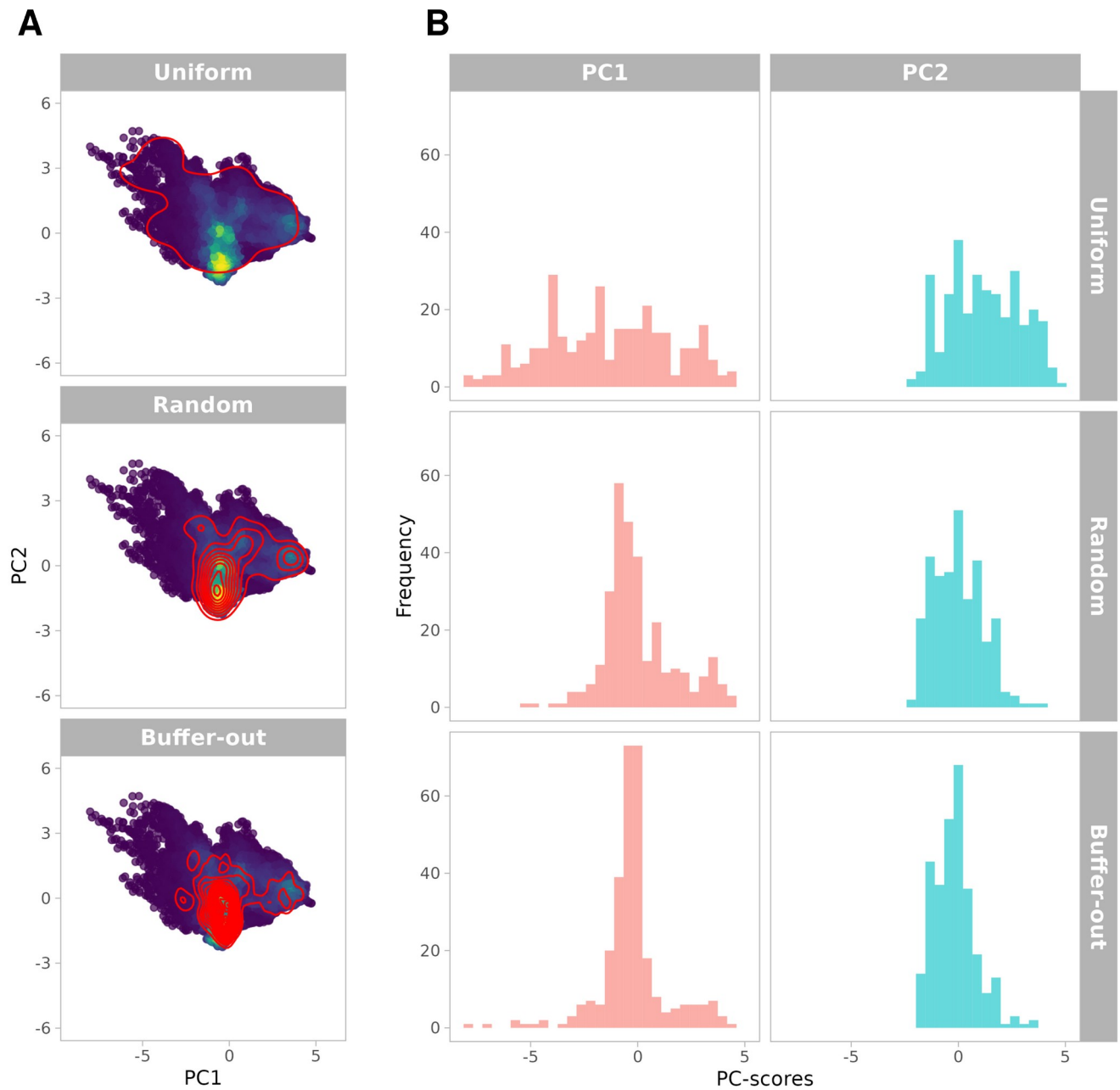
USE: a novel approach to uniformly sample the environmental space

380 3.2 Effect of sample location bias and class overlap

381 The bivariate density plots of the pseudo-absences sampled within the environmental and
382 geographical space highlighted that the uniform approach had the widest and most
383 homogeneous coverage of environmental conditions throughout the environmental space
384 (Fig. 4, see Figure S1.2 for a more detailed representation of the density of pseudo-
385 absences sampled within the environmental space when running the uniform approach). In
386 contrast, the random and buffer-out approaches appeared to be prone to sample location
387 bias, with peaks of high density of pseudo-absences occurring in specific areas of the
388 environmental space, i.e. those associated with the most frequent habitat conditions
389 encountered within the geographical space (Fig. 4).

390 Regarding class overlap, we detected a statistically significant difference in the overlap
391 between the portions of the environmental space occupied by presences and pseudo-
392 absences (one-way ANOVA $F(2, 27) = 5.83$, $p\text{-value} = 0.008$). Specifically, the uniform
393 approach exhibited the lowest overlap in comparison to the other sampling strategies (Fig.
394 5).

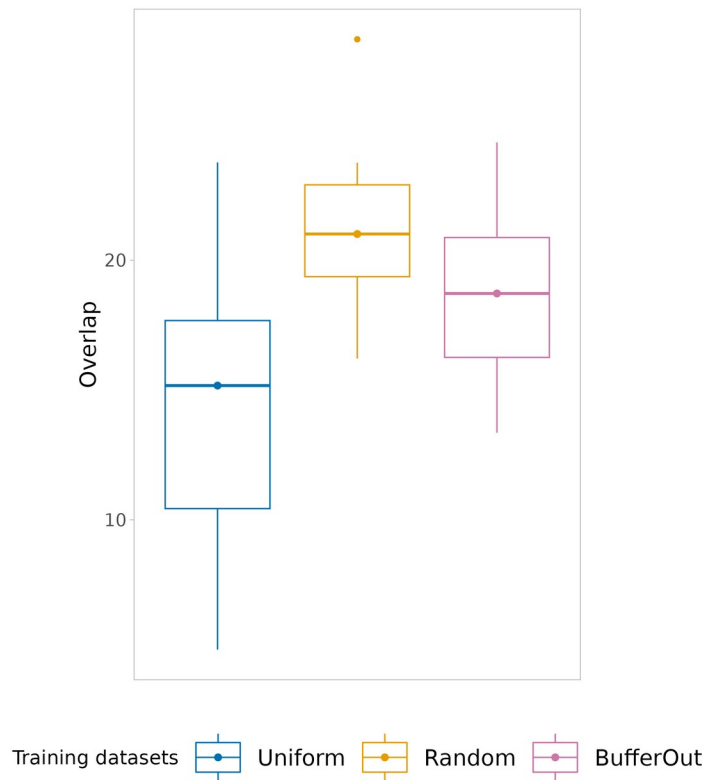
USE: a novel approach to uniformly sample the environmental space



395 **Figure 4:** A) Bivariate plots showing the environmental space generated by a PCA carried out on 5
396 bioclimatic variables. Red lines represent the density of pseudo-absences for an individual VS, as
397 sampled by the random and buffer-out approaches within the geographical space, and by the
398 uniform approach within the environmental space. A more detailed representation of the density of
399 pseudo-absences sampled by the uniform approach is reported in Figure S1.2. B) Histograms

USE: a novel approach to uniformly sample the environmental space

400 showing the frequency distribution of the first two principal components (columns) associated with
401 the different sampling strategies (rows).



402 **Figure 5:** Box plots showing the overlap between environmental spaces generated by
403 presences and pseudo-absences of the VS. Colours are associated with the three sampling
404 strategies used to generate the pseudo-absences (uniform in blue, random in yellow and
405 buffer-out in pink). Dots represent median values of overlap, as computed across 10 VS.

406

USE: a novel approach to uniformly sample the environmental space

407 4 Discussion

408 In this study, we proposed the uniform approach as an alternative strategy to sample
409 pseudo-absences within the environmental space. In contrast to existing techniques, our
410 approach systematically samples pseudo-absences from portions of the environmental
411 space excluding the conditions that are likely to be suitable for the species to establish. As a
412 result, the uniform approach reduces the chance of including false-absences in the training
413 dataset. From a more theoretical perspective, data collected after the application of the
414 kernel-based filter are much closer to the concept of pseudo-absences than those obtained
415 through traditional, geographical sampling approaches. Our findings show that the uniform
416 approach represents a valid strategy for gathering pseudo-absences, resulting in out-of-
417 sample predictive accuracy comparable to the sampling strategies implemented within the
418 geographical space. In addition, the uniform sampling significantly reduces sample location
419 bias and class overlap, which is critical to obtain ecologically meaningful pseudo-absences.
420 Importantly, the uniform approach is flexible, as it allows the user to set parameters (e.g.,
421 kernel bandwidth, sample prevalence, sampling grid resolution) that control how pseudo-
422 absences are sampled within the environmental space. Such flexibility is particularly
423 valuable to mimic different ecological processes that are easier to capture within the
424 environmental space than within the geographical space (e.g., source-sink dynamics). In all
425 cases, by generating informative pseudo-absences, the uniform approach allows satisfying
426 one of the most critical assumptions underpinning habitat suitability modelling: the need for
427 adequate species distribution attributes (i.e., pseudo-absence data here) to model the
428 species-environment relationship (Guisan et al., 2017).

429 4.1 Effect of the sampling approaches on models' predictive performances

USE: a novel approach to uniformly sample the environmental space

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

441 In this sense, results for the CBI, which is currently the go-to accuracy metric for validating
442 HSMs fitted on pseudo-absences (or background points), and for the RMSE were
443 particularly encouraging since the uniform approach scored, together with the buffer-out
444 approach, the highest CBI values and lowest RMSE values across all modelling techniques.
445 The high predictive performance associated with the uniform approach can be attributed to
446 its two main underlying properties: the systematic sampling of the environmental space and
447 the kernel-based filter on the presence data. In particular, the coherent trends of the
448 predictive accuracy metrics (Fig. 3) and the environmental overlap among pseudo-
449 absences collected through the different sampling approaches (Fig. 5) highlight the
450 importance of the kernel-based filter for favouring the discrimination between the
451 environmental features associated with presences and pseudo-absences.

452 Notwithstanding the positive results obtained in terms of predictive performance, we argue
453 that a comparison of metrics of model predictive accuracy may not be the best means for

USE: a novel approach to uniformly sample the environmental space

454 evaluating the adequacy of different sampling strategies carried out within the
455 environmental rather than the geographical space. Indeed, previous studies showed that
456 these metrics are affected by several factors, including sample prevalence (Guisan et al.,
457 2017; Leroy et al., 2018), sample bias (Dubos et al., 2022, Rocchini et al., 2023) or the
458 spatial extent of the study area (Lobo et al., 2008). Moreover, AUC and TSS tend to score
459 high even in case of poor models calibrated on data exhibiting strong sample location bias
460 (Fourcade et al., 2018, Jiménez-Valverde, 2021). Assessing HSMS predictive performance
461 using a set of different predictive accuracy metrics might help the user to critically evaluate
462 the outputs of the models.

463 4.2 Effect of the uniform sampling on sample location bias and class 464 overlap

465 The uniform approach proved to significantly reduce sample location bias, since pseudo-
466 absences were homogeneously scattered across the bivariate density plot of the two
467 principal component axes (Fig. 4a,b, Fig. S1.2 in Supplementary Materials). On the
468 contrary, the two sampling approaches carried out within the geographical space exhibited
469 prominent peaks of density of pseudo-absences in correspondence with the most frequently
470 encountered environmental conditions within the geographical space. As a consequence,
471 the random and buffer-out approaches are likely to provide sub-optimal pseudo-absences
472 for modelling the species-environment relationship (Thuiller et al. 2004; Austin 2007). This
473 aspect gets increasingly relevant as environmental conditions are more heterogeneously
474 distributed across the geographical space (e.g., in mountain regions with high topographic
475 heterogeneity). Therefore, HSMS calibrated on training datasets adequately representing
476 environmental variability rather than wide geographical coverage represent a crucial step to
477 better capture and discriminate species niche breadth (Tessarolo et al., 2014, 2021; Varela

USE: a novel approach to uniformly sample the environmental space

478 et al., 2014; Bazzichetto et al., 2022; Perret and Sax 2022).

479 The uniform approach proved to also significantly reduce class overlap. The `thres`
480 argument passed to the `paSampling` function controls the portion of the environmental
481 space associated with the species presence, thus inherently limiting class overlap by the
482 exclusion of environmental conditions suitable to the species (see Fig. 1c, Fig. 5 and Fig.
483 S1.3). This results in a set of pseudo-absences theoretically much closer to the species'
484 true absences. Given that presence points are unevenly distributed within the environmental
485 space, different kernel thresholds might also be used to handle the sampling of pseudo-
486 absences under particular scenarios. As an example, setting a low kernel threshold would
487 allow excluding accidental presences from unsuitable locations (e.g., 'sink populations')
488 from the training dataset, while potentially including observations from these areas as
489 pseudo-absences. Unfortunately, there is no *a priori* choice about the value of the threshold
490 without having preliminary information on species' ecology, the study area and the goal of
491 the research. For this reason, we provided the `thresh.inspect` function, which produces
492 plots depicting the entire environmental space alongside the portion that would be excluded
493 based on a specific kernel density threshold.

494 4.4 Limitations and usage notes

495 4.4.1 Limitations

496 The first limitation of the uniform approach, which is anyway a general limitation in HSMs
497 (e.g., Cayuela et al., 2012), is that its effectiveness depends on the amount (sample size)
498 and quality (e.g., geographically unbiased data *sensu* Fourcade 2014) of presence data.
499 Indeed, if few presence data are available and/or presence data are geographically biased,
500 the kernel-based filter might not accurately delimit the area associated with suitable

USE: a novel approach to uniformly sample the environmental space

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

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

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

514 Finally, although the idea behind USE and the uniform sampling approach is to provide
515 users with an easy-to-use tool to generate more ecologically meaningful pseudo-absences, we
516 acknowledge the existence of other techniques designed to avoid generating pseudo-absences
517 altogether. Notable examples are point-process analyses (e.g., Isaac et al., 2020), which model
518 the density of presence-only points per unit area, rather than the probability of presences and
519 (pseudo-)absences. More recently, machine-learning methods based on isolation forests were
520 also proposed, with the R package `ITSdm` specifically dedicated to HSMs (Song and Estes,
521 2023). We believe, however, that our approach provides a simpler and more intuitive way to
522 deal with the issue of presence-only data, and thus has a lower threshold for end-users to
523 implement in their workflow.

524 4.4.2 Usage notes

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

USE: a novel approach to uniformly sample the environmental space

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

539 5 Conclusion

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

USE: a novel approach to uniformly sample the environmental space

548 6 Declaration

- 549 • Ethics approval and consent to participate: Not applicable.
- 550 • Competing interests: No conflict of interest has been declared by the authors.
- 551 • Funding: DDR was supported by a FRS-FNRS Belgian grant, ET is supported by the
552 Estonian Research Council grant (MOBJD1030), MB acknowledges funding from the
553 European Union's Horizon Europe research and innovation programme under the Marie
554 Skłodowska-Curie grant agreement No 101066324.
- 555 • Authors' contribution: MB conceived the idea of the uniform approach and wrote the
556 related R functions, while ET and DDR integrated the kernel density-based estimation of
557 presences and the prevalence-related settings. DDR, ET and MB performed the
558 simulations, analysed the data and assembled the `USE` R package. JL, JJJ, SOV, and
559 DR critically commented on the results of the analyses and their interpretation; DDR, ET
560 and MB led the writing of the manuscript and produced a first draft, which was further
561 improved by all other authors.
- 562 • Acknowledgments: The authors are grateful to Joaquin Hortal, who provided
563 constructive feedback and commented on a previous version of this manuscript.
564 Simulations were carried out using the facilities of the High-Performance Computing
565 Center of the University of Tartu.

USE: a novel approach to uniformly sample the environmental space

566 7 Code and Data availability

567 The scripts for replicating the analyses presented in this paper are available at
568 https://github.com/danddr/USE_paper, as well as all the raw outputs of the simulations and
569 statistical analyses (which are available as an .RDS file).

570 We provide a general tutorial to explain how to apply the USE package at
571 https://danddr.github.io/USE/articles/USE_vignette.html. In addition, we provide a tutorial on
572 how to apply the uniform approach based on a real species (the European beech, *Fagus*
573 *sylvatica* L.) in S5. The R script related to the tutorial is available at
574 https://github.com/danddr/USE_paper.

USE: a novel approach to uniformly sample the environmental space

575 References

- 576 Acevedo, P., Jiménez-Valverde, A., Lobo, J. M., and Real, R. (2012). Delimiting the
577 geographical background in species distribution modelling. *Journal of biogeography*,
578 39(8):1383–1390.
- 579 Albert, C. H., Yoccoz, N. G., Edwards Jr, T. C., Graham, C. H., Zimmermann, N. E., and
580 Thuiller, W. (2010). Sampling in ecology and evolution – bridging the gap between theory
581 and practice. *Ecography*, 33(6):1028–1037.
- 582 Austin, M. (2007). Species distribution models and ecological theory: a critical assessment
583 and some possible new approaches. *Ecological Modelling*, 200(1-2), 1-19.
- 584 Barbet-Massin, M., Jiguet, F., Albert, C. H., and Thuiller, W. (2012). Selecting pseudo
585 absences for species distribution models: how, where and how many? *Methods in
586 Ecology and Evolution*, 3(2):327–338.
- 587 Baker, D. J., Maclean, I. M. D., Goodall, M., & Gaston, K. J. (2022). Correlations between
588 spatial sampling biases and environmental niches affect species distribution models.
589 *Global Ecology and Biogeography*, 00, 1– 13.
- 590 Batista, E., Lopes, A., Miranda, P., & Alves, A. (2023). Can species distribution models be
591 used for risk assessment analyses of fungal plant pathogens? A case study with three
592 Botryosphaeriaceae species. *European Journal of Plant Pathology*, 165(1), 41-56.
- 593 Bazzichetto, M., Lenoir, J., Da Re, D., Tordoni, E., Rocchini, D., Malavasi, M., Barták, V. &
594 Sperandii, M. G. (2022). Sampling strategy matters to accurately estimate response
595 curves' parameters in species distribution models. *EcoRxiv*
596 <https://doi.org/10.32942/osf.io/rhys3>
- 597 Bazzichetto, M., Massol, F., Carboni, M., Lenoir, J., Lembrechts, J. J., Joly, R., & Renault,
598 D. (2021). Once upon a time in the far south: Influence of local drivers and functional
599 traits on plant invasion in the harsh sub-Antarctic islands. *Journal of Vegetation
600 Science*, 32(4), e13057.
- 601 Beck, J., Böller, M., Erhardt, A., and Schwanghart, W. (2014). Spatial bias in the GBIF
602 database and its effect on modeling species' geographic distributions. *Ecological
603 Informatics*, 19:10–15.
- 604 Bedia, J., Herrera, S., and Gutiérrez, J. M. (2013). Dangers of using global bioclimatic
605 datasets for ecological niche modeling. limitations for future climate projections. *Global
606 and Planetary Change*, 107:1–12.
- 607 Blonder B, Morrow wcfCB, Harris DJ, Brown S, Butruille G, Laini A, Chen D (2022).
608 _hypervolume: High Dimensional Geometry, Set Operations, Projection, and Inference
609 Using Kernel Density Estimation, Support Vector Machines, and Convex Hulls_. R
610 package version 3.0.4, <<https://CRAN.R-project.org/package=hypervolume>>.
- 611 Booth, T. H., Nix, H. A., Busby, J. R., and Hutchinson, M. F. (2014). Bioclim: the first
612 species distribution modelling package, its early applications and relevance to most

USE: a novel approach to uniformly sample the environmental space

- 613 current maxent studies. *Diversity and Distributions*, 20(1):1–9.
- 614 Cayuela, L., Golicher, D. J., Newton, A. C., Kolb, M., De Albuquerque, F. S., Arets, E. J. M. M.,
615 Alkemade, J. R. M. & Pérez, A. M. (2009). Species distribution modeling in the tropics: problems,
616 potentialities, and the role of biological data for effective species conservation. *Tropical*
617 *Conservation Science*, 2(3), 319-352.
- 618 Cobos, M. E., Peterson, A. T., Barve, N., & Osorio-Olvera, L. (2019). kuenm: an R package for
619 detailed development of ecological niche models using Maxent. *PeerJ*, 7, e6281.
- 620 Da Re, D., Tordoni, E., De Pascalis, F., Negrín-Pérez, Z., Fernández-Palacios, J. M.,
621 Arévalo, J. R., ... & Bacaro, G. (2020). Invasive fountain grass (*Pennisetum setaceum*
622 (Forssk.) Chiov.) increases its potential area of distribution in Tenerife island under future
623 climatic scenarios. *Plant Ecology*, 221(10), 867-882.
- 624 Dinno, A. (2017). *dunn.test: Dunn's Test of Multiple Comparisons Using Rank Sums*. R
625 package version 1.3.5, <https://CRAN.R-project.org/package=dunn.test>.
- 626 Dubos, N., Préau, C., Lenormand, M., Papuga, G., Monsarrat, S., Denelle, P., ... & Luque,
627 S. (2022). Assessing the effect of sample bias correction in species distribution models.
628 *Ecological Indicators*, 145, 109487.
- 629 Duffy, G. A., Coetzee, B. W., Latombe, G., Akerman, A. H., McGeoch, M. A., & Chown, S.
630 L. (2017). Barriers to globally invasive species are weakening across the
631 Antarctic. *Diversity and Distributions*, 23(9), 982-996.
- 632 Duong, T. (2021). *ks: Kernel Smoothing*. R package version 1.13.3.
- 633 Fei, S. and Yu, F. (2016). Quality of presence data determines species distribution model
634 performance: a novel index to evaluate data quality. *Landscape Ecology*, 31(1):31–42.
- 635 Fick, S. E. and Hijmans, R. J. (2017). Worldclim 2: new 1-km spatial resolution climate
636 surfaces for global land areas. *International Journal of Climatology*, 37(12):4302–4315.
- 637 Fourcade, Y. (2021). Fine-tuning niche models matters in invasion ecology. A lesson from
638 the land planarian *Obama nungara*. *Ecological Modelling*, 457, 109686.
- 639 Fourcade, Y., Besnard, A. G., and Secondi, J. (2018). Paintings predict the distribution of
640 species, or the challenge of selecting environmental predictors and evaluation statistics.
641 *Global Ecology and Biogeography*, 27(2):245–256.
- 642 Fourcade, Y., Engler, J. O., Rödder, D., & Secondi, J. (2014). Mapping species distributions
643 with MAXENT using a geographically biased sample of presence data: a performance
644 assessment of methods for correcting sampling bias. *PloS ONE*, 9(5), e97122.
- 645 Grimmett, L., Whitsed, R., & Horta, A. (2020). Presence-only species distribution models
646 are sensitive to sample prevalence: Evaluating models using spatial prediction stability
647 and accuracy metrics. *Ecological Modelling*, 431, 109194.
- 648 Guisan, A., Thuiller, W., and Zimmermann, N. E. (2017). *Habitat suitability and distribution*
649 *models: with applications in R*. Cambridge University Press.
- 650 Hallgren, W., Santana, F., Low-Choy, S., Zhao, Y., and Mackey, B. (2019). Species

USE: a novel approach to uniformly sample the environmental space

- 651 distribution models can be highly sensitive to algorithm configuration. *Ecological*
652 *Modelling*, 408:108719.
- 653 Hattab, T., Garzón-López, C. X., Ewald, M., Skowronek, S., Aerts, R., Horen, H., Brasseur,
654 B., Gallet-Moron, E., Spicher, F., Decocq, G., et al. (2017). A unified framework to model
655 the potential and realized distributions of invasive species within the invaded range.
656 *Diversity and Distributions*, 23(7):806–819.
- 657 Hortal, J., Jiménez-Valverde, A., Gómez, J. F., Lobo, J. M., and Baselga, A. (2008).
658 Historical bias in biodiversity inventories affects the observed environmental niche of the
659 species. *Oikos*, 117(6):847–858.
- 660 Hysen, L., Nayeri, D., Cushman, S., & Wan, H. Y. (2022). Background sampling for multi-scale
661 ensemble habitat selection modeling: Does the number of points matter?. *Ecological Informatics*,
662 72, 101914.
- 663 Isaac, N. J., Jarzyna, M. A., Keil, P., Dambly, L. I., Boersch-Supan, P. H., Browning, E., ... &
664 O'Hara, R. B. (2020). Data integration for large-scale models of species distributions. *Trends in*
665 *ecology & evolution*, 35(1), 56-67.
- 666 Iturbide, M., Bedia, J., Herrera, S., del Hierro, O., Pinto, M., and Gutiérrez, J. M. (2015). A
667 framework for species distribution modelling with improved pseudo-absence generation.
668 *Ecological Modelling*, 312:166–174.
- 669 Jackson, S. T. and Overpeck, J. T. (2000). Responses of plant populations and
670 communities to environmental changes of the late quaternary. *Paleobiology*, 26(S4):194–
671 220.
- 672 Jarvie, S., & Svenning, J. C. (2018). Using species distribution modelling to determine
673 opportunities for trophic rewilding under future scenarios of climate change. *Philosophical*
674 *Transactions of the Royal Society B: Biological Sciences*, 373(1761), 20170446.
- 675 Jiménez-Valverde, A. (2021). Prevalence affects the evaluation of discrimination capacity in
676 presence-absence species distribution models. *Biodiversity and Conservation*, 30(5),
677 1331-1340.
- 678 Jiménez-Valverde, A., Lobo, J. M., & Hortal, J. (2008). Not as good as they seem: the
679 importance of concepts in species distribution modelling. *Diversity and*
680 *Distributions*, 14(6), 885-890.
- 681 Jiménez-Valverde, A., Acevedo, P., Barbosa, A. M., Lobo, J. M., and Real, R. (2013).
682 Discrimination capacity in species distribution models depends on the representativeness
683 of the environmental domain. *Global Ecology and Biogeography*, 22(4):508–516.
- 684 Kass, J. M., Muscarella, R., Galante, P. J., Bohl, C. L., Pinilla-Buitrago, G. E., Boria, R. A., Soley-
685 Guardia, M., & Anderson, R. P. (2021). ENMeval 2.0: Redesigned for customizable and
686 reproducible modeling of species' niches and distributions. *Methods in Ecology and*
687 *Evolution*, 12(9), 1602-1608.
- 688 Leroy, B., Delsol, R., Hugueny, B., Meynard, C. N., Barhoumi, C., Barbet-Massin, M., and
689 Bellard, C. (2018). Without quality presence–absence data, discrimination metrics such
690 as tss can be misleading measures of model performance. *Journal of Biogeography*,

USE: a novel approach to uniformly sample the environmental space

- 691 45(9):1994–2002.
- 692 Leroy, B., Meynard, C. N., Bellard, C., and Courchamp, F. (2016). virtualspecies, an r
693 package to generate virtual species distributions. *Ecography*, 39(6):599–607.
- 694 Lobo, J. M., Jiménez-Valverde, A., and Hortal, J. (2010). The uncertain nature of absences
695 and their importance in species distribution modelling. *Ecography*, 33(1):103–114.
- 696 Lobo, J. M., Jiménez-Valverde, A., and Real, R. (2008). Auc: a misleading measure of the
697 performance of predictive distribution models. *Global Ecology and Biogeography*,
698 17(2):145– 151.
- 699 Magri, D., Vendramin, G. G., Comps, B., Dupanloup, I., Geburek, T., Gömöry, D., ... & De
700 Beaulieu, J. L. (2006). A new scenario for the Quaternary history of European beech
701 populations: palaeobotanical evidence and genetic consequences. *New Phytologist*,
702 171(1), 199-221.
- 703 Mammola, S. and Cardoso, P. (2020). Functional diversity metrics using kernel density n-
704 dimensional hypervolumes. *Methods in Ecology and Evolution*, 11(8):986–995.
- 705 Meynard, C. N., Leroy, B., and Kaplan, D. M. (2019). Testing methods in species
706 distribution modelling using virtual species: what have we learnt and what are we
707 missing? *Ecography*, 42(12):2021–2036.
- 708 Naimi, B. and Araújo, M. B. (2016). sdm: a reproducible and extensible r platform for
709 species distribution modelling. *Ecography*, 39(4):368–375.
- 710 Newbold, T. (2018). Future effects of climate and land-use change on terrestrial vertebrate
711 community diversity under different scenarios. *Proceedings of the Royal Society B*,
712 285(1881):20180792.
- 713 Perret, D. L. and Sax, D. F. (2022). Evaluating alternative study designs for optimal
714 sampling of species' climatic niches. *Ecography*.
- 715 Poli et al. (2022) Coupling fossil records and traditional discrimination metrics to test how
716 genetic information improves species distribution models of the European beech *Fagus*
717 *sylvatica*. *European Journal of Forest Research*, 141: 253–265
- 718 Phillips, S. J., Anderson, R. P., Dudík, M., Schapire, R. E., and Blair, M. E. (2017). Opening
719 the black box: An open-source release of maxent. *Ecography*, 40(7):887–893.
- 720 Phillips, S. J., Dudík, M., Elith, J., Graham, C. H., Lehmann, A., Leathwick, J., and Ferrier,
721 S. (2009). Sample selection bias and presence-only distribution models: implications for
722 background and pseudo-absence data. *Ecological Applications*, 19(1):181–197.
- 723 R Core Team (2021). *R: A Language and Environment for Statistical Computing*. R
724 Foundation for Statistical Computing, Vienna, Austria.
- 725 Rocchini, D., Tordoni, E., Marchetto, E. *et al.* A quixotic view of spatial bias in modelling the
726 distribution of species and their diversity. *npj biodiversity* 2, 10 (2023).
727 <https://doi.org/10.1038/s44185-023-00014-6>
- 728 Ronquillo, C., Alves-Martins, F., Mazimpaka, V., Sobral-Souza, T., Vilela-Silva, B., Medina,

USE: a novel approach to uniformly sample the environmental space

- 729 N. G., and Hortal, J. (2020). Assessing spatial and temporal biases and gaps in the
730 publicly available distributional information of iberian mosses. *Biodiversity Data Journal*,
731 8.
- 732 Santini, L., Benítez-López, A., Maiorano, L., Čengić, M., and Huijbregts, M. A. (2021).
733 Assessing the reliability of species distribution projections in climate change research.
734 *Diversity and Distributions*, 27(6):1035–1050.
- 735 Scott, D.W. (1992) *Multivariate Density Estimation: Theory, Practice, and Visualization*,
736 John Wiley & Sons
- 737 Sillero, N. and Barbosa, A. M. (2020). Common mistakes in ecological niche models.
738 *International Journal of Geographical Information Science*, pages 1–14.
- 739 Song, L., & Estes, L. (2023). ITSDM: Isolation forest-based presence-only species
740 distribution modelling and explanation in R. *Methods in Ecology and Evolution*, 14(3),
741 831-840.
- 742 Støa, B., Halvorsen, R., Stokland, J. N., and Gusarov, V. I. (2019). How much is enough?
743 influence of number of presence observations on the performance of species distribution
744 models. *Sommerfeltia*, 39(1):1–28.
- 745 Svenning, J.-C. and Skov, F. (2004). Limited filling of the potential range in European tree
746 species. *Ecology Letters*, 7(7):565–573.
- 747 Tessarolo, G., Lobo, J. M., Rangel, T. F., and Hortal, J. (2021). High uncertainty in the
748 effects of data characteristics on the performance of species distribution models.
749 *Ecological Indicators*, 121:107147.
- 750 Tessarolo, G., Rangel, T. F., Araújo, M. B., and Hortal, J. (2014). Uncertainty associated
751 with survey design in species distribution models. *Diversity and Distributions*,
752 20(11):1258–1269.
- 753 Thuiller, W., Brotons, L., Araújo, M. B., & Lavorel, S. (2004). Effects of restricting
754 environmental range of data to project current and future species distributions.
755 *Ecography*, 27(2), 165– 172. <https://doi.org/10.1111/j.0906-7590.2004.03673.x>
- 756 Valavi, R., Elith, J., Lahoz-Monfort, J. J., and Guillera-Arroita, G. (2021). Modelling species
757 presence-only data with random forests. *Ecography*, 44(12):1731–1742.
- 758 VanDerWal, J., Shoo, L. P., Graham, C., & Williams, S. E. (2009). Selecting pseudo-
759 absence data for presence-only distribution modeling: how far should you stray from what
760 you know?. *Ecological Modelling*, 220(4), 589-594.
- 761 Varela, S., Anderson, R. P., García-Valdés, R., and Fernández-González, F. (2014).
762 Environmental filters reduce the effects of sampling bias and improve predictions of
763 ecological niche models. *Ecography*, 37(11):1084–1091.
- 764 Venables WN, Ripley BD (2002). *Modern Applied Statistics with S*, Fourth edition. Springer,
765 New York. ISBN 0-387-95457-0, <https://www.stats.ox.ac.uk/pub/MASS4/>.
- 766 Wasof et al. (2015) Disjunct populations of European vascular plant species keep the same

USE: a novel approach to uniformly sample the environmental space

767 climatic niches. *Global Ecology and Biogeography*, 24: 1401-1412

768 Wickham H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. ISBN
769 978-3-319-24277-4, <https://ggplot2.tidyverse.org>.

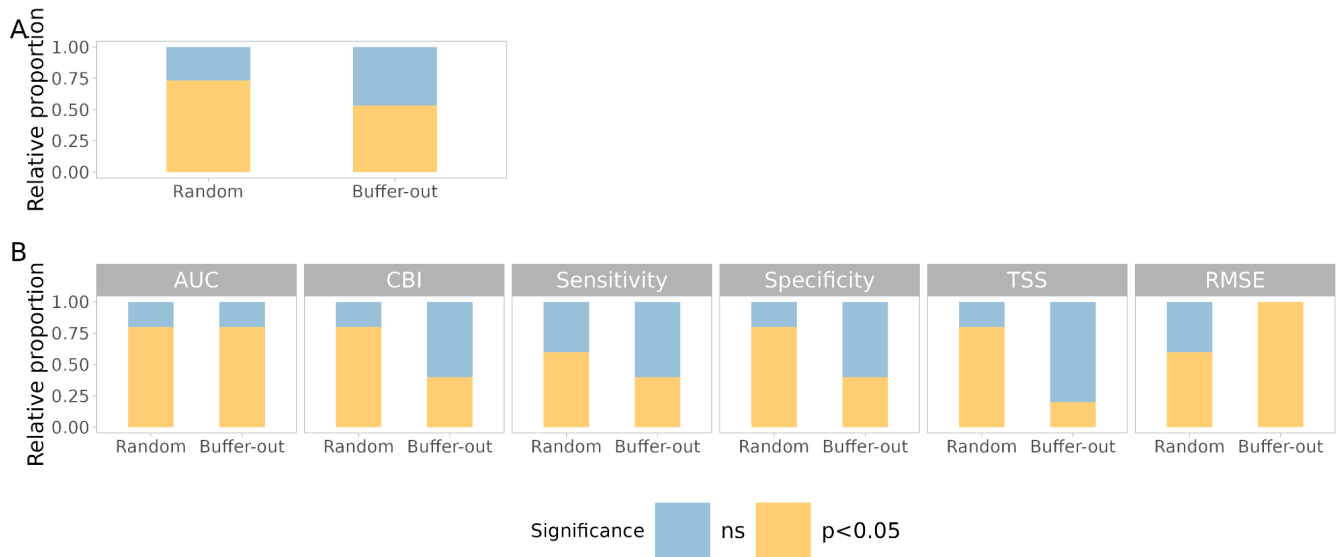
Supplementary Material 1

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 with the other different sampling strategies. P-values were adjusted using Holm correction.

Model	Metric	Comparisons	χ^2	P.val
BRT	AUC	BufferOut - Uniform	93.89	ns
BRT	AUC	Random - Uniform	93.89	p<0.05
BRT	CBI	BufferOut - Uniform	46.45	ns
BRT	CBI	Random - Uniform	46.45	p<0.05
BRT	RMSE	BufferOut - Uniform	6.08	p<0.05
BRT	RMSE	Random - Uniform	6.08	ns
BRT	Sensitivity	BufferOut - Uniform	1.97	ns
BRT	Sensitivity	Random - Uniform	1.97	ns
BRT	Specificity	BufferOut - Uniform	91.85	p<0.05
BRT	Specificity	Random - Uniform	91.85	p<0.05
BRT	TSS	BufferOut - Uniform	87.43	ns
BRT	TSS	Random - Uniform	87.43	p<0.05
GAM	AUC	BufferOut - Uniform	101.78	p<0.05
GAM	AUC	Random - Uniform	101.78	p<0.05
GAM	CBI	BufferOut - Uniform	29.26	p<0.05
GAM	CBI	Random - Uniform	29.26	p<0.05
GAM	RMSE	BufferOut - Uniform	25.43	p<0.05
GAM	RMSE	Random - Uniform	25.43	ns
GAM	Sensitivity	BufferOut - Uniform	28.82	p<0.05
GAM	Sensitivity	Random - Uniform	28.82	p<0.05
GAM	Specificity	BufferOut - Uniform	101.68	ns
GAM	Specificity	Random - Uniform	101.68	p<0.05
GAM	TSS	BufferOut - Uniform	91.55	ns

Model	Metric	Comparisons	χ^2	P.val
GAM	TSS	Random - Uniform	91.55	p<0.05
GLM	AUC	BufferOut - Uniform	27.47	p<0.05
GLM	AUC	Random - Uniform	27.47	ns
GLM	CBI	BufferOut - Uniform	44.97	ns
GLM	CBI	Random - Uniform	44.97	p<0.05
GLM	RMSE	BufferOut - Uniform	35.25	p<0.05
GLM	RMSE	Random - Uniform	35.25	p<0.05
GLM	Sensitivity	BufferOut - Uniform	8.18	ns
GLM	Sensitivity	Random - Uniform	8.18	p<0.05
GLM	Specificity	BufferOut - Uniform	53.85	p<0.05
GLM	Specificity	Random - Uniform	53.85	ns
GLM	TSS	BufferOut - Uniform	25.54	p<0.05
GLM	TSS	Random - Uniform	25.54	ns
Maxent	AUC	BufferOut - Uniform	101.78	p<0.05
Maxent	AUC	Random - Uniform	101.78	p<0.05
Maxent	CBI	BufferOut - Uniform	90.22	p<0.05
Maxent	CBI	Random - Uniform	90.22	ns
Maxent	RMSE	BufferOut - Uniform	31.96	p<0.05
Maxent	RMSE	Random - Uniform	31.96	p<0.05
Maxent	Sensitivity	BufferOut - Uniform	24.26	p<0.05
Maxent	Sensitivity	Random - Uniform	24.26	p<0.05
Maxent	Specificity	BufferOut - Uniform	88.04	ns
Maxent	Specificity	Random - Uniform	88.04	p<0.05
Maxent	TSS	BufferOut - Uniform	89.73	ns
Maxent	TSS	Random - Uniform	89.73	p<0.05
RF	AUC	BufferOut - Uniform	100.83	p<0.05
RF	AUC	Random - Uniform	100.83	p<0.05
RF	CBI	BufferOut - Uniform	93.99	ns

Model	Metric	Comparisons	χ^2	P.val
RF	CBI	Random - Uniform	93.99	p<0.05
RF	RMSE	BufferOut - Uniform	31.28	p<0.05
RF	RMSE	Random - Uniform	31.28	p<0.05
RF	Sensitivity	BufferOut - Uniform	0.21	ns
RF	Sensitivity	Random - Uniform	0.21	ns
RF	Specificity	BufferOut - Uniform	97.83	ns
RF	Specificity	Random - Uniform	97.83	p<0.05
RF	TSS	BufferOut - Uniform	87.93	ns
RF	TSS	Random - Uniform	87.93	p<0.05



S1.1: 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 with the other different 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.

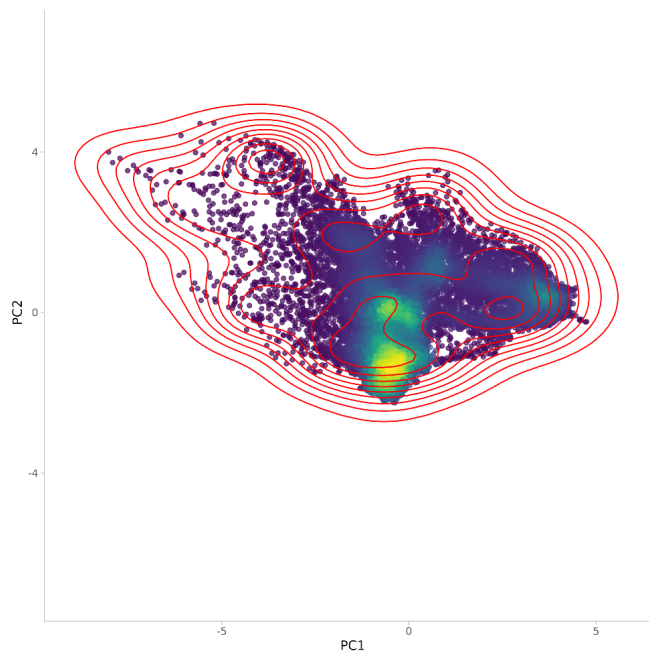


Figure S1.2: Virtual Species PC-scores bivariate plot for the Uniform approach only.

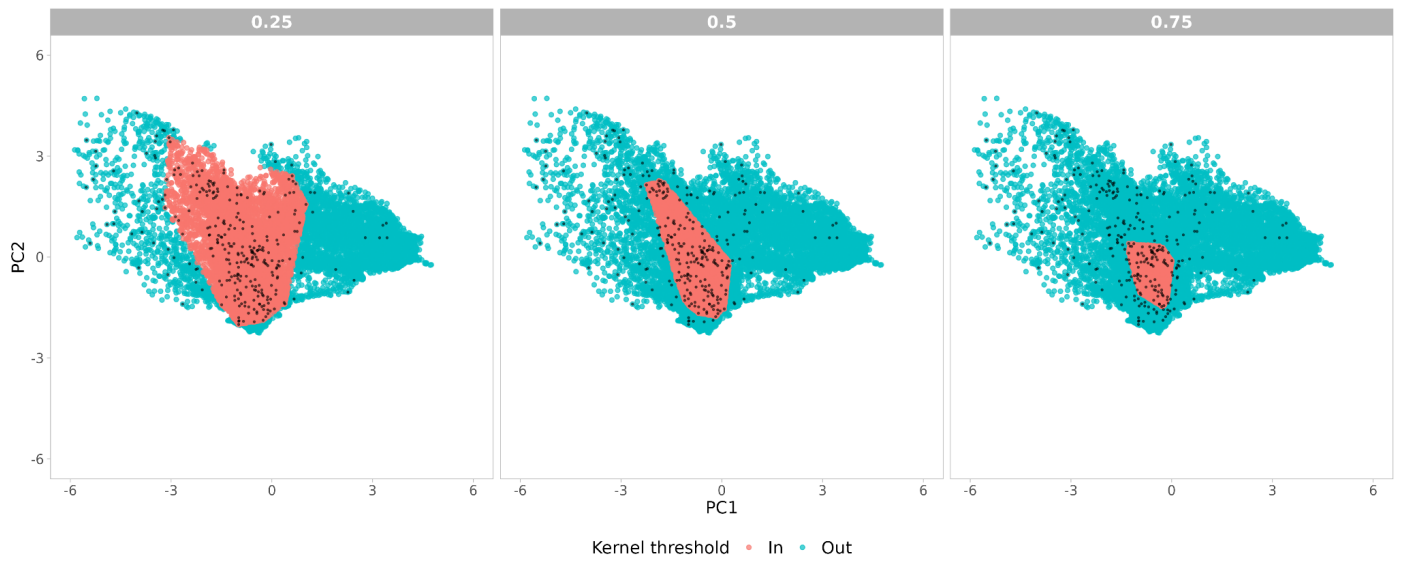


Figure S1.3: Effect of setting different kernel thresholds on the inclusion/exclusion of pseudo-absences eventually sampled using the uniform approach (black dots are the true VS presences displayed within the environmental space). Setting a low value of the kernel threshold (e.g., 0.25) increases the portion of the environmental space excluded from the uniform sampling, while, on the contrary, 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, we also repeated the entire workflow on 10 VS with two different prevalence values. Specifically, in both cases we kept a training dataset consisting of 300 presences but we used alternatively 600 and 3,000 pseudo-absences (sample prevalence = 0.5 and sample prevalence = 0.1, respectively).

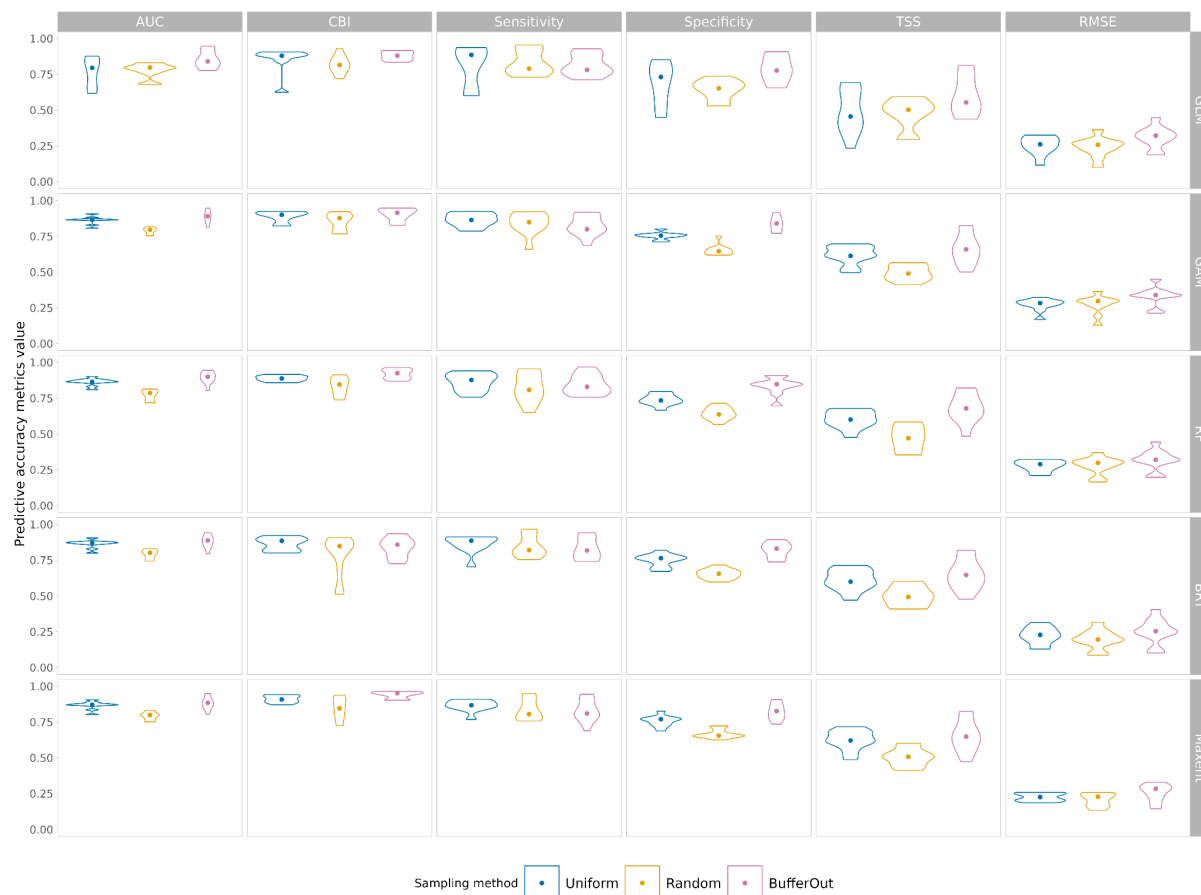


Figure S2.1: Violin plots reporting the distribution of the values of the metrics of predictive performance for the HSMs of the 10 VS (the dots are the median values), considering 5 predictors, and using a sample prevalence = 0.5. Columns indicate the different performance metrics while rows report the modelling techniques used to compute HSMs.

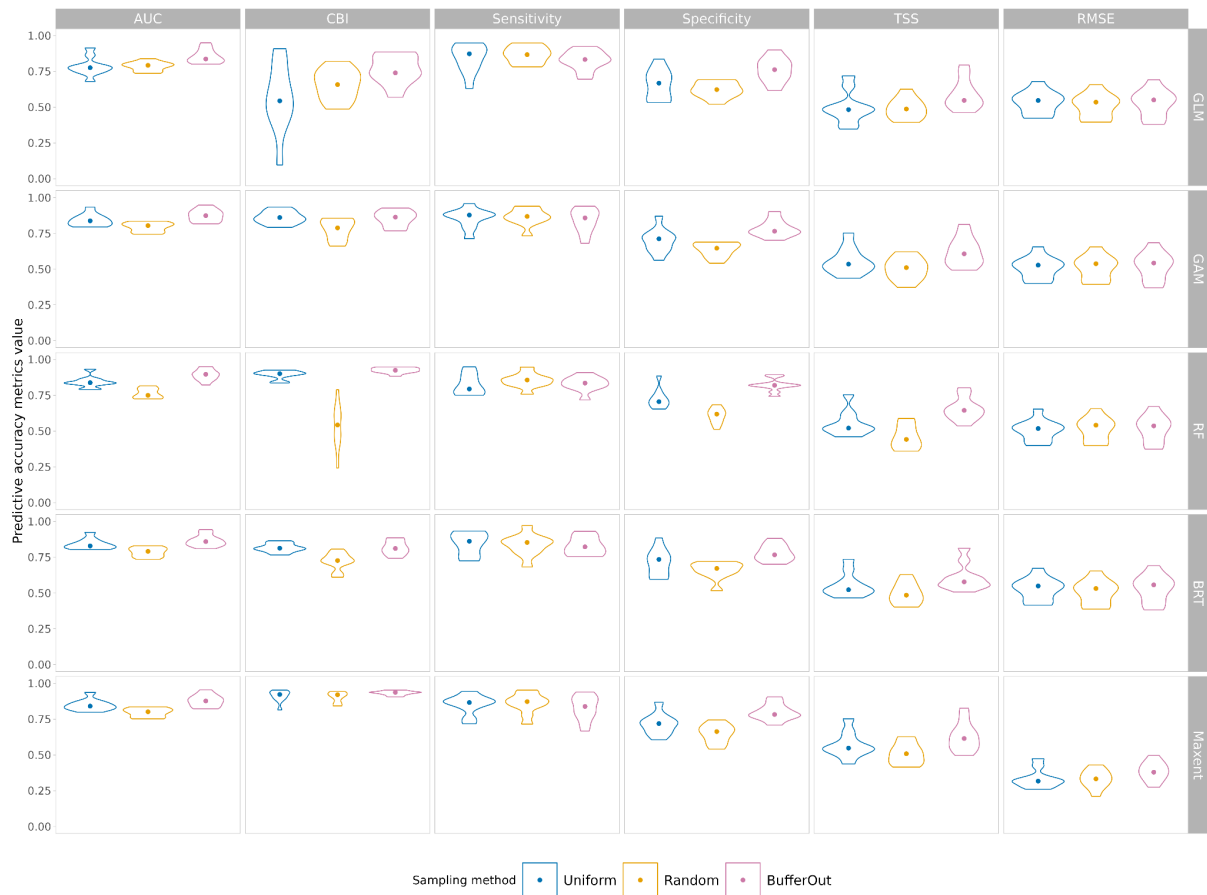


Figure S2.2: Violin plots reporting the distribution of the values of the metrics of predictive performance for the HSMs of the 10 VS (the dots are the median values), considering 5 predictors, and using a sample prevalence = 0.1. Columns indicate the different performance metrics while rows report the modelling techniques used to compute HSMs.

Supplementary Material 3

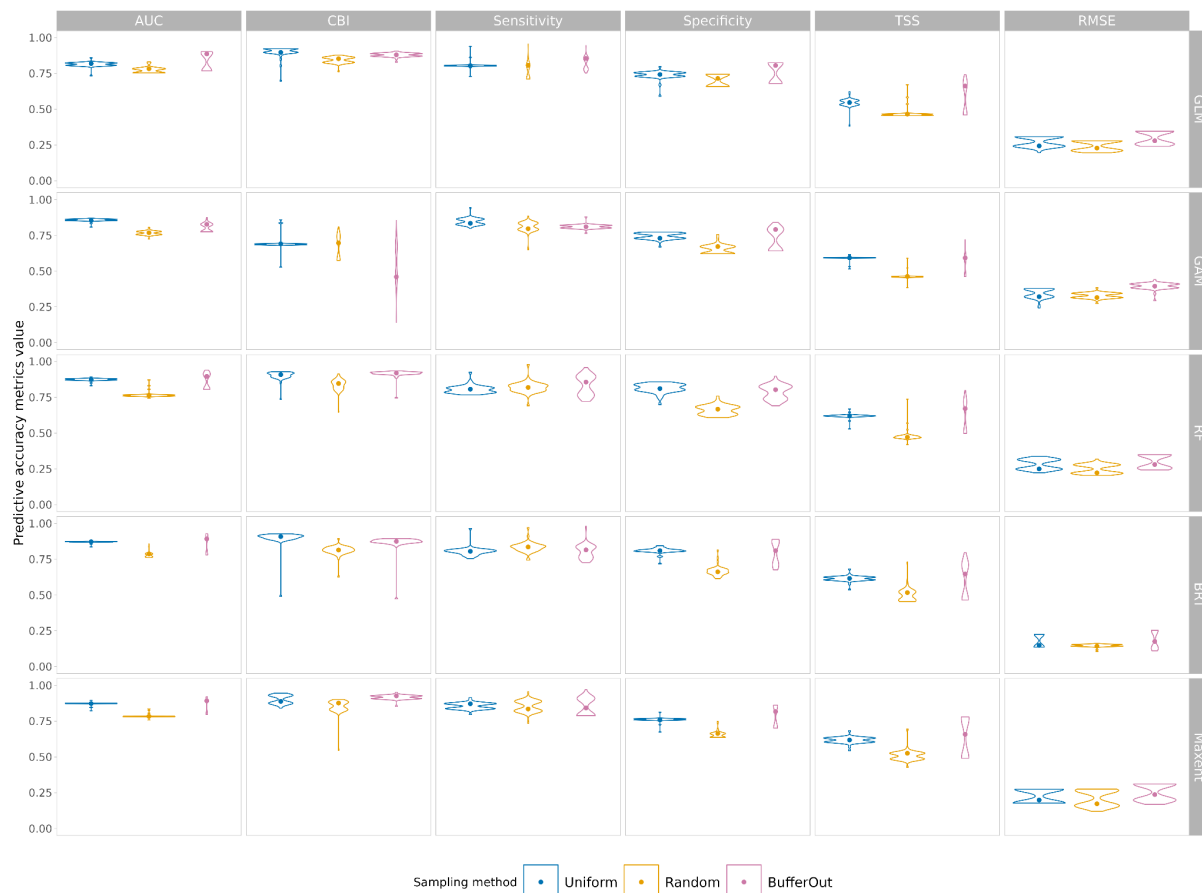


Fig. S3.1: Violin plots reporting the distribution of the values of the metrics of predictive performance for the HSMs of the 50 VS modelled as a function of the 19 bioclimatic predictors, and setting sample prevalence equal to 1 (i.e., same number of presences and pseudo-absences). Dots represent median values of the metrics of predictive accuracy. Columns indicate the different performance metrics, while rows the modelling techniques used to compute HSMs. AUC = Area Under the Curve; CBI = Continuous Boyce Index, TSS = True Skill Statistic; RMSE = Root Mean Squared Error; GLM = Generalized Linear Model; GAM = Generalized Additive Model; RF = Random Forest; BRT = Boosted Regression Trees.

Supplementary Material 4

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

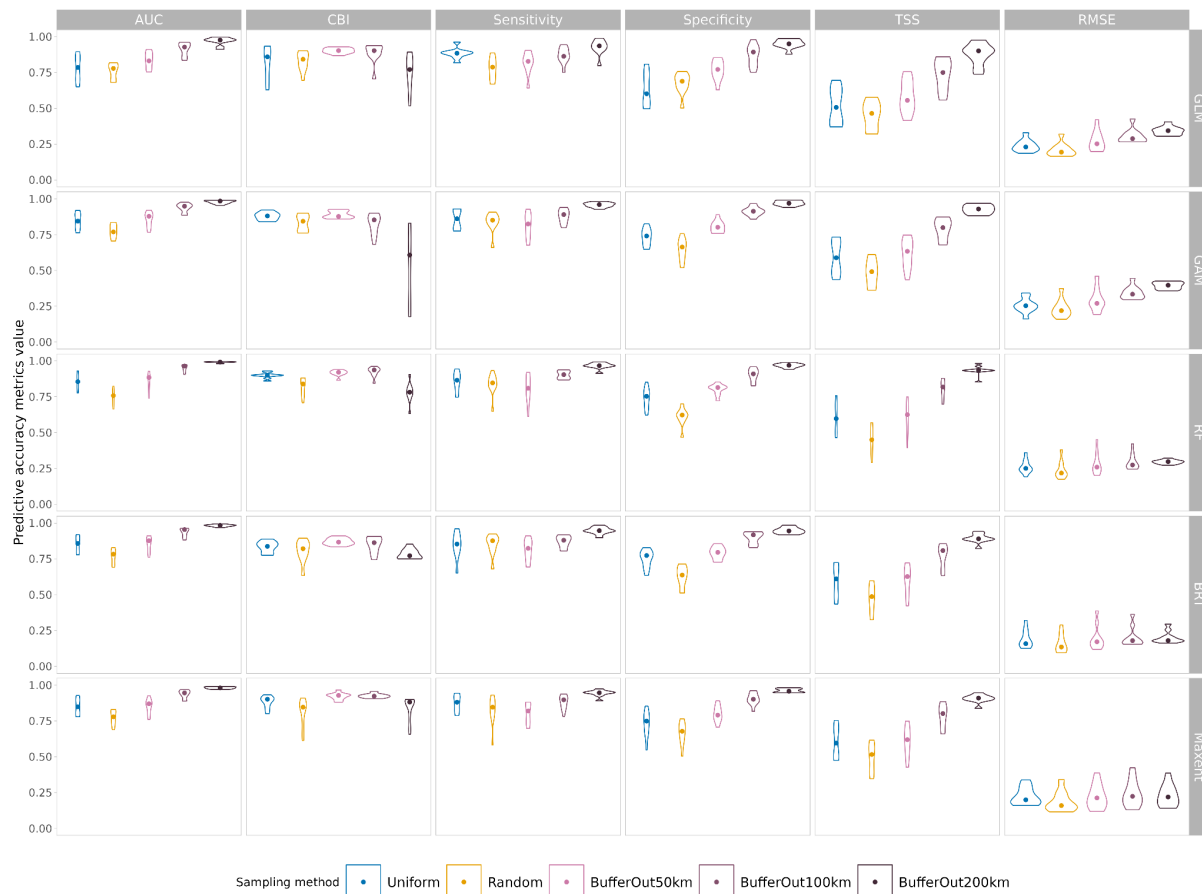


Fig. S4.1: Violin plots reporting the distribution of the values of the metrics of predictive performance for the HSMs of the 10 VS modelled as a function of 5 bioclimatic predictors, and setting sample prevalence equal to 1 (i.e., same number of presences and pseudo-absences). We varied the size of the radius for the buffer-out approach, namely: 50, 100 and 200 km. Dots represent median values of the metrics of predictive accuracy. Columns indicate the different performance metrics, while rows the modelling techniques 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 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 1×1 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.

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 HSMs for *Fagus sylvatica* (GLM and RF). Models' predictive performance was assessed through internal (5-fold CV 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	

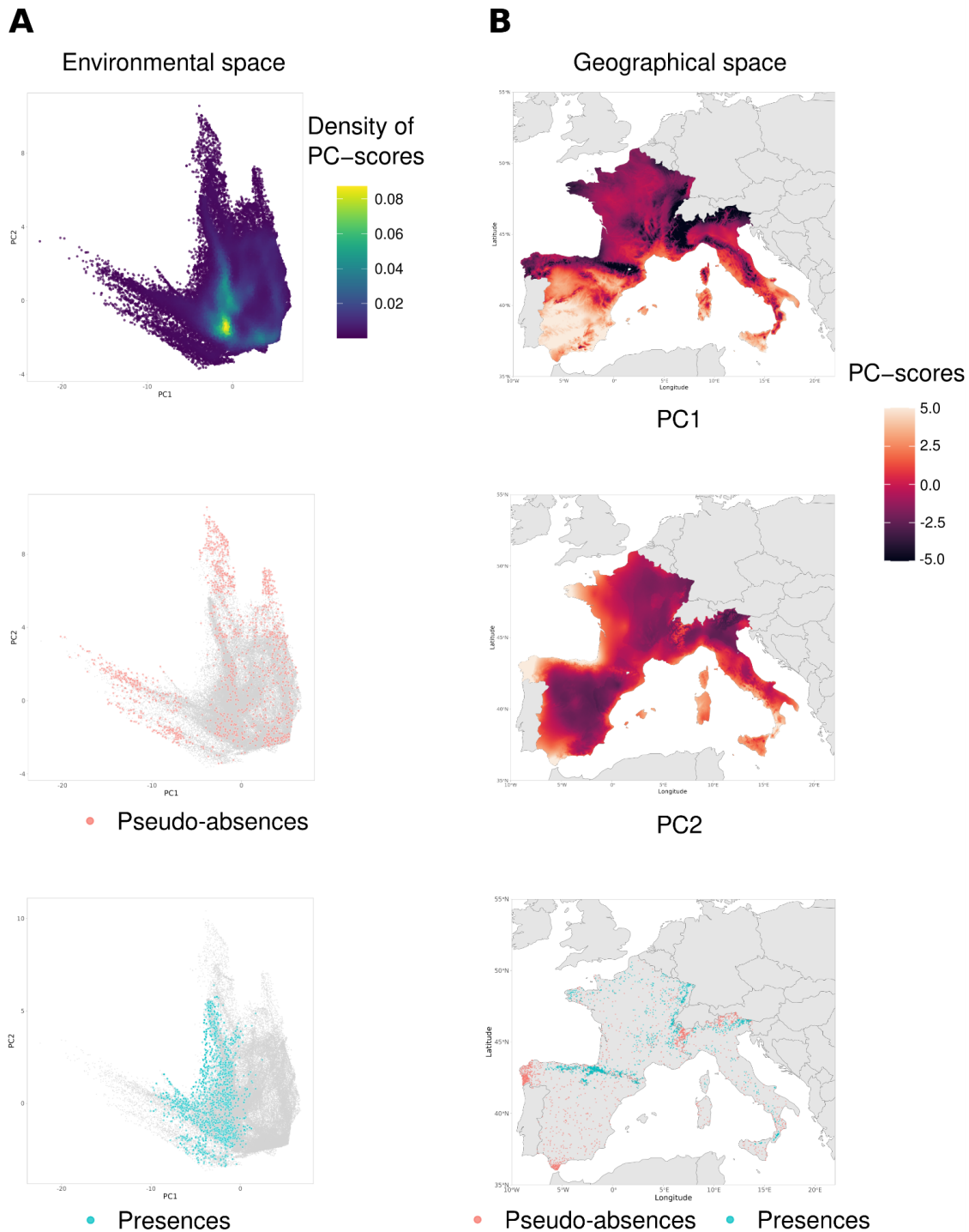


Figure S5.2: (A) environmental space available for *Fagus sylvatica* in Italy, Spain and France, and the position of presences (light blue) and pseudo-absences (red) sampled within the environmental space using the uniform approach; (B) distribution of PC-scores across the geographical space, and location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