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 Lenoir[3],

5 Jonas J. Lembrechts[4], Sophie O. Vanwambeke[1,],

6 Duccio Rocchini[5,6], and Manuele Bazzichetto[7,8]†

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] Department of Ecology and Global Change, Centro de Investigaciones sobre
19 Desertificacion, Moncada (Valencia), Spain
20 [8] Faculty of Environmental Sciences, Department of Spatial Sciences, Czech University of
21 Life Sciences Prague, Kamýcka 129, 16500, Praha-Suchdol, Czech Republic
22
23 †DDR, ET and MB equally contributed to this work
24
25 **Corresponding author**: Daniele Da Re, daniele.dare@uclouvain.be
26
27
28

29

# Abstract

1. Correlative habitat suitability models infer the geographical distribution of species using occurrence data and environmental variables. While data on species presence are increasingly accessible, the difficulty to confirm real absences in the field often forces researchers to generate them *in silico*. To this aim, pseudo-absences are commonly randomly sampled across the study area (i.e., the geographical space). However, this introduces sample location bias (i.e., the sampling is unbalanced towards the most frequent habitats occurring within the geographical space) and class overlap (i.e., overlap between environmental conditions associated with species presences and pseudo-absences) in the training dataset.

2. To mitigate this, we propose an alternative methodology (i.e., uniform approach) that systematically samples pseudo-absences within a portion of the environmental space delimited by a kernel-based filter, which minimises the number of false-absences included in the training set.

3. We simulated 50 virtual species and modelled their distribution using training datasets assembled with the occurrences of the virtual species and pseudo-absences collected using the uniform approach and other approaches that randomly sample pseudo-absences within the geographical space. We compared the predictive performance of the models and evaluated the extent of sample location bias and class overlap associated with the different sampling strategies.

4. Results indicated that the uniform approach: (i) effectively reduces sample location bias and class overlap; (ii) provides comparable predictive accuracy than sampling

53        strategies carried out in the geographic space; (iii) ensures gathering pseudo-

54        absences adequately representing the environmental conditions available across the

55        study area. We developed a set of R functions in an accompanying R package called

56        `USE` to disseminate the uniform approach.

57    **Keywords**: background points, pseudo-absence, ecological niche models,

58    environmental space, habitat suitability models, presence-only models, sample

59    location bias, class overlap, species distribution models, reproducibility.

60 # 1 Introduction

61 Correlative habitat suitability models (hereafter, HSMs) are a class of statistical models

62 used to describe the relationship between species attributes (e.g., presence-absence,

63 abundance) and a set of spatially-explicit variables chiefly representing abiotic and human-

64 related factors (e.g., climate, soil, land-use). These models are rooted in the niche theory

65 (i.e., *Hutchinsonian* niche, see Guisan et al., 2017) and rely on both theoretical and

66 practical assumptions: (i) species are assumed to be at (quasi)equilibrium with their

67 environment (Hattab et al., 2017); (ii) the set of predictors used to fit HSMs includes all

68 necessary information to capture the ecological niche of the species; and (iii) species

69 attributes, used as the response variable, need to be appropriate for the intended model

70 purpose (e.g., biodiversity conservation, forecasting biological invasions, assessing the

71 effects of global change; Tessarolo et al., 2021; see also Guisan et al., 2017 for a thorough

72 review on the theoretical assumptions underpinning HSMs). Some of these assumptions

73 are hardly, if ever, met in nature since species are seldom at equilibrium with their

74 environment (Svenning and Skov, 2004), posing several limitations to the use and

75 interpretation of HSMs' outputs. Acknowledging and, when possible, addressing these

76 limitations still makes HSMs a powerful toolbox for understanding the drivers of the species'

77 realized and potential distributions (*sensu* Jackson and Overpeck, 2000). For this reason,

78 HSMs are still widely applied in several research fields, including biogeography (Wasof et

79 al., 2015; Duffy et al., 2017), climate change ecology (Jarvie and Svenning, 2018),

80 conservation biology (Newbold, 2018; Santini et al., 2021), and invasion ecology (Hattab et

81 al., 2017; Da Re et al. 2020; Bazzichetto et al. 2021).

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

83 biological data for modelling the ecological niche of the species, which means that species

84  attributes, being either presence-absence or abundance data, should allow effectively

85  describing the true species-environment relationship (Guisan et al., 2017; Baker et al.,

86  2022). However, while information on species occurrence (i.e., presence) is usually readily

87  accessible through field-collected observations or museum/herbaria records, trustworthy

88  absence data are by far more difficult to gather or to confirm in the field (Jiménez-Valverde

89  et al., 2008), as their sampling requires labour-intensive and costly field campaigns (Hattab

90  et al., 2017). The usual lack of true absence data has led to the development of HSMs

91  approaches that either rely solely on presence data (so-called 'presence-only models', such

92  as the BIOCLIM model; Booth et al. 2014) or combine presence data with pseudo-absences

93  or background points for modelling species distributions (e.g., the MaxEnt algorithm; Phillips

94  et al., 2017). The terms pseudo-absences and background points are often used as

95  synonyms in the scientific literature (Sillero and Barbosa, 2020), yet these two concepts

96  reflect rather different conditions. On the one hand, pseudo-absences are sampled from

97  geographical locations that are thought to feature unsuitable environmental conditions for

98  the species to establish (Barbet-Massin et al., 2012). On the other hand, background points

99  are collected from the whole spectrum of environmental conditions present in the study

100  area, thereby possibly including suitable locations for the species (i.e., presence locations;

101  Phillips et al., 2009; Hallgren et al., 2019). Therefore, the use of pseudo-absences rather

102  than background points reflects the user's degree of uncertainty about the species'

103  ecological preferences, with  background points being used when there is no *a priori*

104  knowledge about the unsuitable environmental conditions for the species. Although we

105  acknowledge the difference between pseudo-absences and background points, for the sake

106  of simplicity and because we feel the concept of pseudo-absence adheres more to what we

107  propose in this study, hereafter we will  always refer to pseudo-absences.

108        To date, the most common approaches for sampling pseudo-absences involve (i)

109   surveying a large sample of points randomly located across the study area (e.g., 10,000;

110   Barbet-Massin et al. 2012; Iturbide et al., 2015; Støa et al., 2019) or sampling them either

111   (ii) within or (iii) outside the area covered by buffers built around presence locations

112   (VanDer Wal et al., 2009; Bedia et al., 2013). Beyond the pros and cons of each individual

113   approach, a common thread is that they all randomly deploy pseudo-absences across the

114   geographical space, which usually results in the oversampling of the most common habitat

115   conditions, namely those that are more geographically widespread throughout the area

116   under investigation (Tessarolo et al., 2014, 2021; Ronquillo et al., 2020). This phenomenon,

117   which is generally known as sample location bias (Elith et al. 2011), has detrimental effects

118   on HSMs for different reasons. First, it determines the incomplete sampling of the

119   environmental conditions actually experienced by a species (i.e., the realised environment

120   *sensu* Jackson and Overpeck, 2000), possibly leading to the estimation of truncated

121   species response curves (Hortal et al., 2008; Albert et al., 2010; Beck et al., 2014). Second,

122   it affects the predictive performance of HSMs (Acevedo et al., 2012), which is reflected in

123   the behaviour of the metrics used to evaluate them (Jiménez-Valverde et al., 2013; Sillero

124   and Barbosa, 2020).

125      To overcome this issue, previous studies (Varela et al. 2014; Hattab et al., 2017)

126   proposed to sample species presence and (true) absence data through a systematic

127   sampling of the environmental conditions available throughout the study area, thus limiting

128   the artificial constraint imposed by the random sampling towards the most widespread

129   environments. More specifically, Varela et al. (2014), Hattab et al. (2017) and Perret and

130   Sax (2022) suggested collecting species' presence and/or absence within 2- or 3-

131   dimensional environmental spaces obtained using ordination techniques. Such approaches

132   significantly contributed to the improvement and standardisation of the way species

133   observations, including pseudo-absences, can be collected to calibrate HSMs reducing

134  sample location bias. Yet, they do not explicitly consider class overlap, another relevant

135  methodological issue encountered when collecting pseudo-absences through random

136  sampling across the geographical space. Class overlap refers to the overlap between

137  environmental conditions associated with species presence and absence, thus hindering

138  the concept of pseudo-absences itself. It has negative effects on the predictive performance

139  of HSMs and is particularly critical for machine learning techniques, while regression

140  techniques such as GLMs seem to be less affected (Barbet-Massin et al., 2012; Grimmett,

141  Whitsed and Horta, 2020; Valavi et al., 2021). So far, class overlap has been addressed

142  using resampling techniques more oriented to adjusting an unbalanced number of classes

143  in the response variable (i.e., the 'up-' or 'down-sampling' approach;  Valavi et al., 2021),

144  irrespective of the technique to obtain pseudo-absences.

145      As far as we know, there are no approaches for sampling pseudo-absences that are

146  able to mitigate both sample location bias and class overlap. Thus, here we present an

147  alternative sampling strategy, which we called the 'uniform' approach, that builds upon

148  existing strategies for systematically sampling the environmental space to select pseudo-

149  absences. The novel aspect of the uniform approach is that, beyond reducing sample

150  location bias, it also minimises class overlap by implementing a kernel-based filter that is

151  used to delineate the portion of the environmental space where to collect pseudo-absences.

152  To test our approach, we simulated 50 virtual species and compared the predictive

153  performance of HSMs trained on pseudo-absences sampled using the uniform approach as

154  well as other sampling strategies traditionally carried out within the geographical space: (i)

155  random (i.e., pseudo-absences randomly sampled within the geographical space); (ii)

156  buffer-in and iii) buffer-out (i.e., pseudo-absences randomly collected within or outside

157  buffers built around presence locations, respectively). To foster reproducibility, we provide

158  an accompanying R package called USE (Uniform Sampling of the Environmental space),

159  which bundles the R functions needed to implement the uniform approach. The package is

160  available at https://github.com/danddr/USE. Finally, we provide a tutorial to explain how to

161  apply the uniform approach to real case studies, using the European beech *Fagus sylvatica*

162  L. as a target species.

# 163  2 Methods

## 164  2.1 Simulation of virtual species

165  We used virtual species (hereafter VS) as they provide the great advantage of knowing the

166  true generative process underlying the species geographical distribution (Meynard et al.,

167  2019). The realised environmental space (*sensu* Jackson and Overpeck 2000) of the 50

168  virtual species was created using the bioclimatic variables gathered from the WorldClim

169  database (www.worldclim.org; spatial resolution ~18.6 km at the Equator; Fick and Hijmans,

170  2017). We restricted the distribution of the simulated VS (and those of the climatic

171  variables) to the geographical extent spanning from -12° W to 25° E and from 36° to 60° N

172  (approximately Western and Southern Europe) to significantly reduce the computational

173  effort to process the entire workflow. Each VS was generated using a random set of five

174  climatic variables (out of the 19) through the function `generateRandomSp` from the R

175  package `virtualspecies` (Leroy et al., 2016), which randomly assigns relationships

176  between the VS and those climatic variables (e.g., linear, quadratic relationships). This way,

177  we obtained a raster layer reporting the habitat suitability index (HSI, Fig. 1a), which we

178  then converted to a binary (i.e., presence-absence) map using  the function `convertToPA`.

179  Further details about parameters setting can be found in the R code available at

180  https://github.com/danddr/USE_paper.

## 181  2.2 Sampling of the pseudo-absences

182  Regardless of the sampling approach and modelling technique used to calibrate the HSMs,

183  the ratio between the number of presences and pseudo-absences in the training datasets

184  (i.e., sample prevalence) was kept equal to 1, which means an equal number of presences

185  and pseudo-absences were collected. In practice, each of the VS-specific training dataset

186  included 300 presences, which were randomly sampled within the geographical extent

187  using the function `sampleOccurrences` from the `virtualspecies` R package.

188  Consequently, we collected an equal number of pseudo-absences according to the four

189  sampling strategies presented below.

190  *2.2.1 Uniform approach: pseudo-absences sampled within the environmental space*

191  For each VS (i.e., iteration), we built a 2-dimensional environmental space by keeping the

192  first two axes of a principal component analysis (PCA) performed on the correlation matrix

193  of the five randomly selected bioclimatic variables used to generate the realised

194  environment (Fig. 1b). Each time, we checked that the first two principal component axes

195  accounted for at least 70% of the total bioclimatic variability. Then, we uniformly sampled

196  pseudo-absences, here intended as the PC-scores projected onto the environmental space,

197  using the function `uniformSampling`. More specifically, each pair of PC-scores

198  represents the position of a given geographical location, as defined by the bioclimatic

199  conditions it features, within the environmental space. Below, we present a step-by-step

200  description of the uniform sampling performed by the function `paSampling`, which

201  internally calls `uniformSampling`, in the `USE` R package:

202      1. First, kernel density estimation is used to calculate the probability density function of

203     the presence data within the 2-dimensional environmental space (Fig. 1c). Similar

204     uses of kernel density estimation have become popular in recent years, especially due

205     to their increasing use in trait-based ecology to compute probabilistic hypervolumes

206     and trait probability densities (Mammola and Cardoso, 2020 and reference therein).

207     All pseudo-absences associated with a probability threshold equal to or greater than

208     0.75 (i.e., the default threshold value used in the `paSampling` function) are excluded,

209     since these points are likely to bear environmental conditions associated with

210     presence locations and can therefore introduce false-absences in the training dataset.

211     The kernel bandwidth (i.e., the width of the kernel density function that defines its

212     shape) can be either defined by the user or automatically estimated by the function

213     `paSampling`. In the latter case, the function uses a bandwidth selector by internally

214     calling the function `Hpi` of the R package `ks` (Duong, 2021).

215   2. A sampling grid constituted by a pre-selected resolution (e.g., 10 X 10 cells) is overlaid

216     on the 2-dimensional environmental space (Fig. 1d). The optimal resolution of the

217     sampling grid can be found using the function `optimRes` from the USE package.

218     This function operates as follows:

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

220       between the pseudo-absences (PC-scores) in the cell and the centroid of their

221       convex hull is computed;

222   -   The same measure is computed in each cell of the sampling grid and the average of

223       the cell-specific averages is computed (hereafter, grid average);

224   -   The procedure above is separately repeated on sampling grids of increasing

225       resolution (i.e., increasing number of cells);

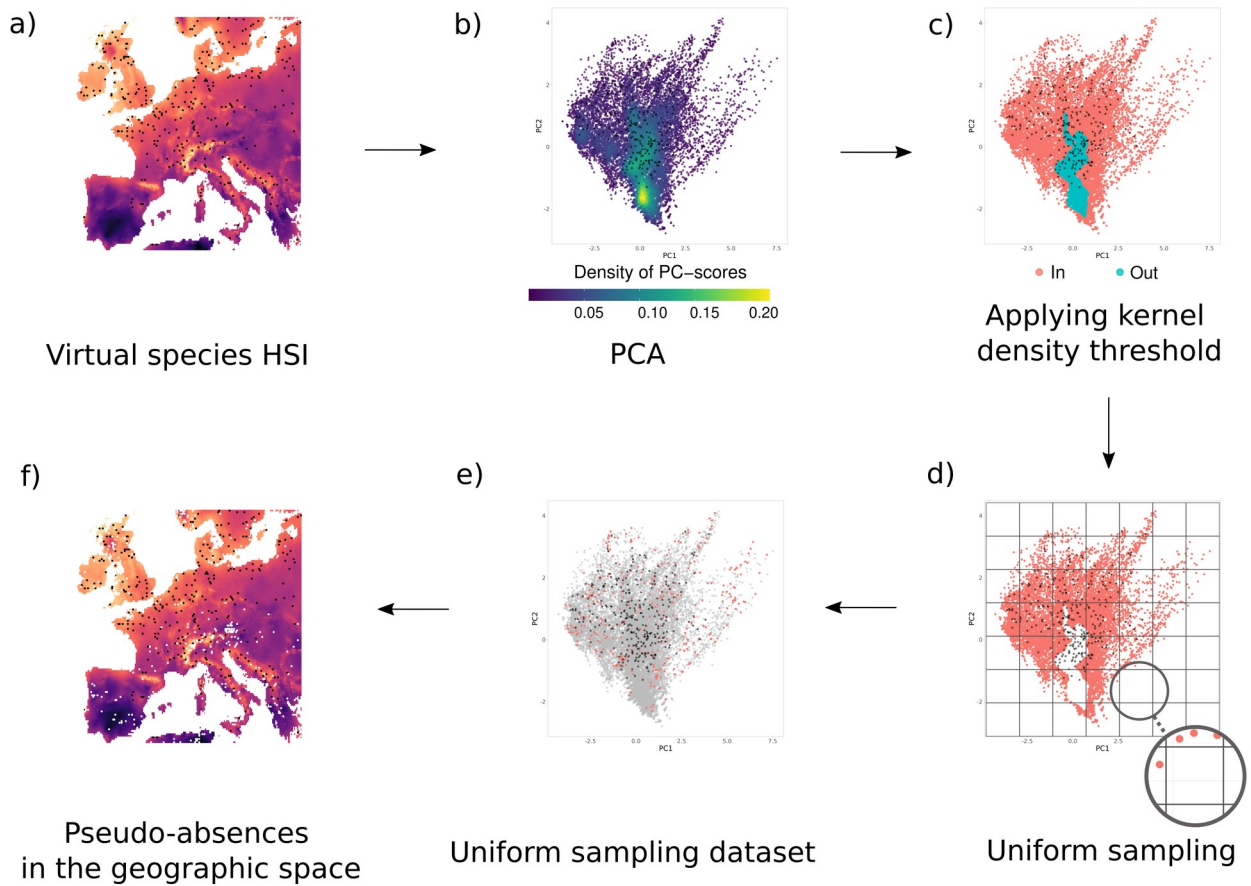226   -   The resulting grid averages are used as a measure of the aggregation among

227    pseudo-absences within the cells of the sampling grids. This value is compared

228    across resolutions and the best grid is chosen as the one providing the best trade-off

229    between resolution and average distance among points within cells (i.e., resolution

230    that allows uniformly sampling the environmental space without overfitting it). More

231    specifically, the best grid is the one whose resolution is just below that which would

232    not allow the average distance among pseudo-absences to be reduced by more than

233    10% (other values can be set by the user).

234   3. Once the resolution is set, the sampling grid is sequentially scanned (i.e., cell by cell)

235    by the `uniformSampling` function called via `paSampling` function and, from each

236    grid cell, a given number of pseudo-absences is randomly collected. At this stage, the

237    pseudo-absences associated with environmental conditions too close to those of the

238    presence locations are already excluded (see step 1). Notice that the pseudo-

239    absences are randomly selected within the area of each cell of the sampling grid, and

240    not at the centroid nor at the nodes.

241   The total number of pseudo-absences sampled within each cell of the sampling grid can be

242   set by the user (using the argument `n.tr`, default `n.tr` = 5), who can also indicate a

243   desired sample prevalence. If the sample prevalence is not specified, fewer pseudo-absences

244   are likely to be eventually sampled than expected (i.e., `n.tr` × number of cells). This

245   behaviour happens because no points are collected in empty cells, and less points than n.tr are

246   available within the cells at the boundary of the environmental space (Figure 1d). Similarly, no

247   pseudo-absences are collected within the core area of the presences (excluded in step 1). If a

248   sample prevalence is set by the user, the sampling grid is surveyed until the set sample

249   prevalence is achieved.

250



Virtual species HSI      PCA      Applying kernel density threshold

Pseudo-absences in the geographic space      Uniform sampling dataset      Uniform sampling

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

252    approach: a) habitat suitability index of the *i-th* virtual species (lighter colours indicate higher

253    habitat suitability and black dots represent presence points in the geographical space); b)

254    PCA performed on the environmental variables in the study region, lighter colours indicate

255    high PC-scores densities and black dots represent the presence points in the environmental

256    space; c) application of the kernel-based filter, which splits the environmental space in sub-

257    spaces associated with either the environmental conditions more suitable for the species

258    (blue) or those associated with less/not suitable environmental conditions (red; black dots

259    represent presence points); d) pseudo-absences are uniformly sampled across a sampling

260  grid overlaid to the 2-dimensional environmental space. Specifically, pseudo-absences are

261  sampled within each cell of the 2-d grid. The inset map shows an example of an empty cell

262  (i.e., a grid cell containing no pseudo-absences; black dots represent presence points); e)

263  the red dots represent the dataset of pseudo-absences collected within the environmental

264  space using the uniform approach;  f) the white dots represent the pseudo-absences

265  collected within the environmental space using the uniform approach displayed in the

266  geographical space, black dots represent VS presence points.


267  *2.2.2 Pseudo-absences sampled within the geographical extent*


268  The sampling of pseudo-absences within the geographical extent defined above was

269  conducted using the random, buffer-in and buffer-out approaches. For the random approach

270  (Barbet-Massin et al. 2012; Iturbide et al., 2015; Støa et al., 2019), we simply located 300

271  random pseudo-absences across the geographical extent. For the other approaches, we

272  created a buffer of 50 km radius around each presence location, and then we randomly

273  sampled the pseudo-absences within (cf. buffer-in; VanDer Wal et al., 2009) and outside (cf.

274  buffer-out; Bedia et al., 2013) the buffers. Notice that for the buffer-out approach pseudo-

275  absences were collected within the convex hull of the species distribution (i.e., the convex

276  hull that connects the outer occurrences of the species and thus delimits the range actually

277  covered by the species in the geographical space).


278  ## 2.3 Comparison among sampling strategies

279  *2.3.1 Predictive performance comparison*

280  The overall workflow of the analyses is described in Fig. 2. For each of the 50 VS and for

281  each of the four sampling strategies (i.e., uniform, random, buffer-in, buffer-out), we built a

282  specific dataset combining the presence records with the pseudo-absences sampled within

283  the environmental and the geographical space (Fig. 1e). First, we modelled the presence–

284  pseudo-absences data as a function of the same five bioclimatic variables used to generate

285  each of the 50 VS. To this aim, we randomly partitioned each dataset (specific for a

286  sampling strategy) in 5 training (70% observations) and testing (30%) sets, which we used

287  to calibrate and validate five modelling algorithms: (i) generalised linear models (GLMs); (ii)

288  generalised additive models (GAMs); (iii) random forests (RFs); (iv) boosted regression

289  trees (BRTs); and (v) MaxEnt. In total, we fitted 5,000 HSMs (50 VS species × 4 different

290  sets of pseudo-absences × 5 modelling algorithms × 5 replicates of 70-30% partitions). To

291  fit the HSMs, we used the R package `sdm` (Naimi and Araújo, 2016). Although we

292  acknowledge the importance of fine-tuning HSMs (Fourcade, 2021), we leave model

293  settings at their default value since it would have been unfeasible to individually

294  parameterise each algorithm for all 50 VS and sampling strategies.

295      Then, we compared the predictive performance of each combination of sampling

296  approaches and modelling techniques computing the following metrics: (i) area under the

297  receiver operating characteristic curve (AUC); (ii) continuous Boyce index (CBI); (iii)

298  sensitivity; (iv) specificity; and (v) true skill statistics (TSS). A detailed description of the five

299  modelling techniques and five validation metrics can be found in Guisan et al. (2017). To

300  compare the predictive performance of the HSMs fitted under different combinations of

301  sampling strategy and modelling technique, we visually assessed the results of the 50 VS

302  simulations using violin plots reporting the distribution of the values of the predictive

303  performance metrics listed above. Furthermore, we tested differences among the predictive

304  performance of the sampling strategies using Kruskall-Wallis tests, followed by Dunn's post

305  hoc rank sum comparisons using the `dunn.test` R package (Dinno, 2017) and correcting

306  p-values for multiple comparisons with the Holm correction.

307      To test the potential effect of varying sample prevalences on our comparison, we

308    repeated the entire workflow on 10 VS using two different prevalence values: 0.5 and 0.1.

309    Specifically, for each VS, we generated two training datasets with 300 presences, but we

310    combined them with 600 and 3,000 pseudo-absences, respectively.

311    *2.3.2 Sample location bias and class overlap*

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

313    strategies, we extracted the pseudo-absences of a single VS and map their spatial

314    aggregation within the environmental space using bivariate density plots. The aim was to

315    identify which, among the four sampling strategies, was more subject to oversampling

316    particular environmental conditions within the geographical space. In principle, the sampling

317    strategies more affected by sample location bias would exhibit a clear aggregation of

318    pseudo-absences within the environmental space. We visually assessed the areas of the

319    environmental space sampled by the different sampling strategies using the function

320    `geom_density_2d` of the `ggplot2` R package (Wickham, 2016). This function performs a

321    2D kernel density estimation using the `kde2d` function of the `MASS` R package (Venables

322    and Ripley, 2002) and displays the results with contours.
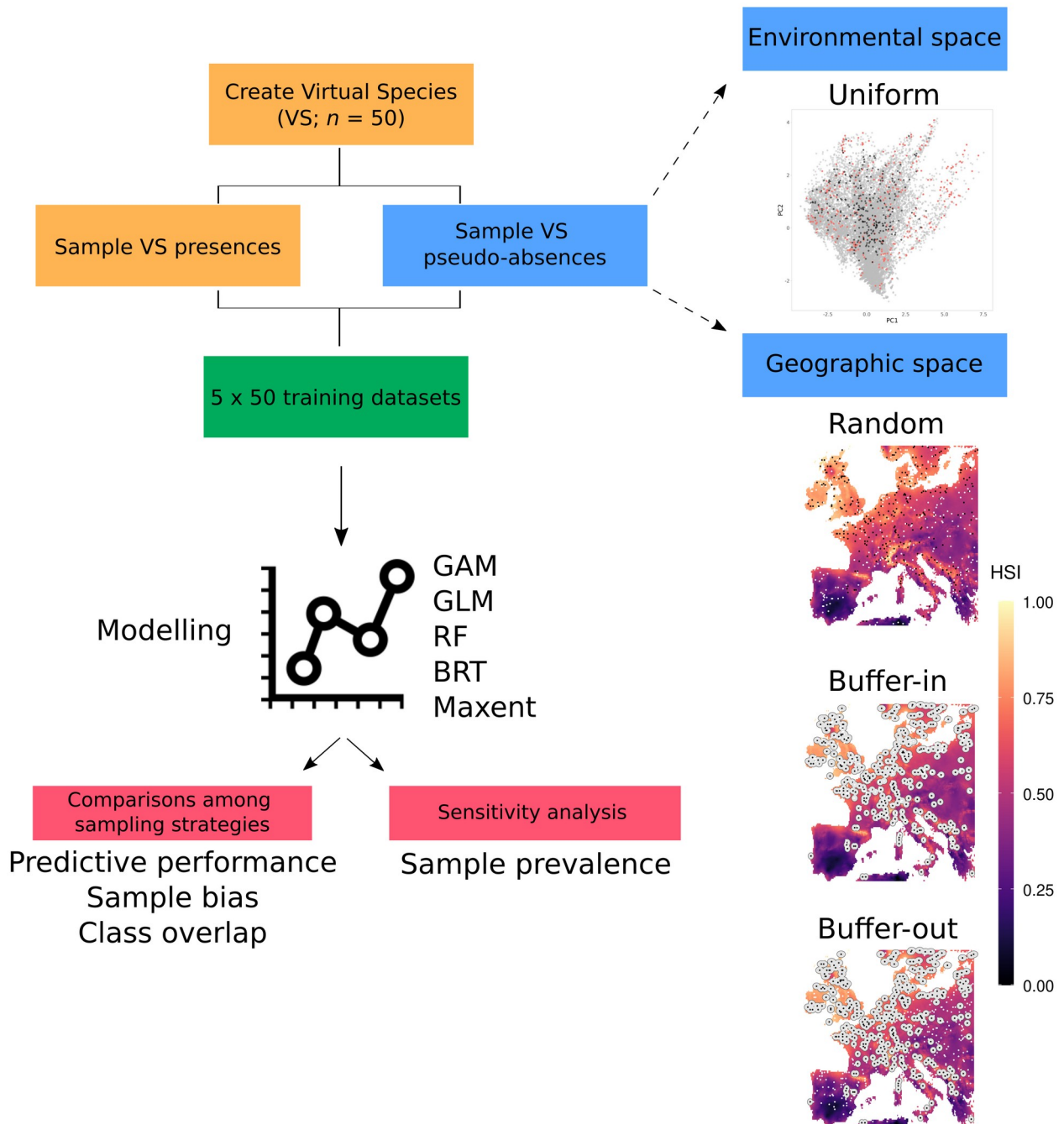
323         To assess the effectiveness of the uniform approach for mitigating class overlap, we

324    simulated 10 further VS, sampled their presences and pseudo-absences using the four

325    sampling strategies and mapped the position of the presence and pseudo-absence points

326    within the environmental space following the procedure explained in section 2.2.1 and figure

327    1a,b. Then, we computed the hypervolume of the presences and pseudo-absences using

328    the `hypervolumes` R package (Blonder, 2022) and calculated the overlap between them.

329    Significant differences in the degree of overlap were tested using one-way ANOVA and

330    Tukey HSD test.


331    ## 2.4 Real-case study

332 To illustrate how to apply the uniform approach with the USE R package, we modelled the

333 realised distribution of *Fagus sylvatica* in Italy, France and Spain (hereafter, Western

334 Europe). We chose *F. sylvatica* as a target species because its distribution and

335 biogeographic history is well-known across Europe (Magri et al., 2006; Poli et al., 2022).

336 The whole procedure is described in S4.

USE: a novel approach to uniformly sampling the environmental space



**Figure 2** Overall workflow of the analysis described in the Methods section.

# 338  3 Results

## 339  3.1 Comparison of the predictive performance associated with geographical vs

## 340  environmental sampling

341  Overall, the uniform approach performed equal to or better than the geographical
342  approaches in terms of out-of-sample prediction. In particular, the uniform and buffer-out
343  strategies showed, on average, the highest predictive accuracy and their performances
344  were not significantly different (Fig. 3). Pairwise comparisons between the performance of
345  the random and buffer-out approaches against the uniform approach showed statistically
346  significant differences in 92% and 72% of the cases (5,000 models obtained from 50 VS x 4
347  sampling strategies x 5 algorithms x 5 replicates), respectively. However, these differences
348  were algorithm- and metric- dependent and did not point to a higher predictive performance
349  of the uniform approach (Fig. 3). The buffer-out and uniform approaches exhibited the most
350  similar values for AUC, sensitivity and TSS, while CBI values tended to be higher with
351  respect to random sampling (i.e., Dunn's test: p-value > 0.05; Tab. S1, Fig. S1.1).  Finally,
352  the buffer-in consistently showed the lowest performance in all comparisons regardless of
353  the algorithm and predictive performance metric used. The observed pattern of the
354  difference among predictive performances was consistent across sample prevalences (Fig.
355  S3.1-3.2).

**Figure 3**: 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 5, randomly selected 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.
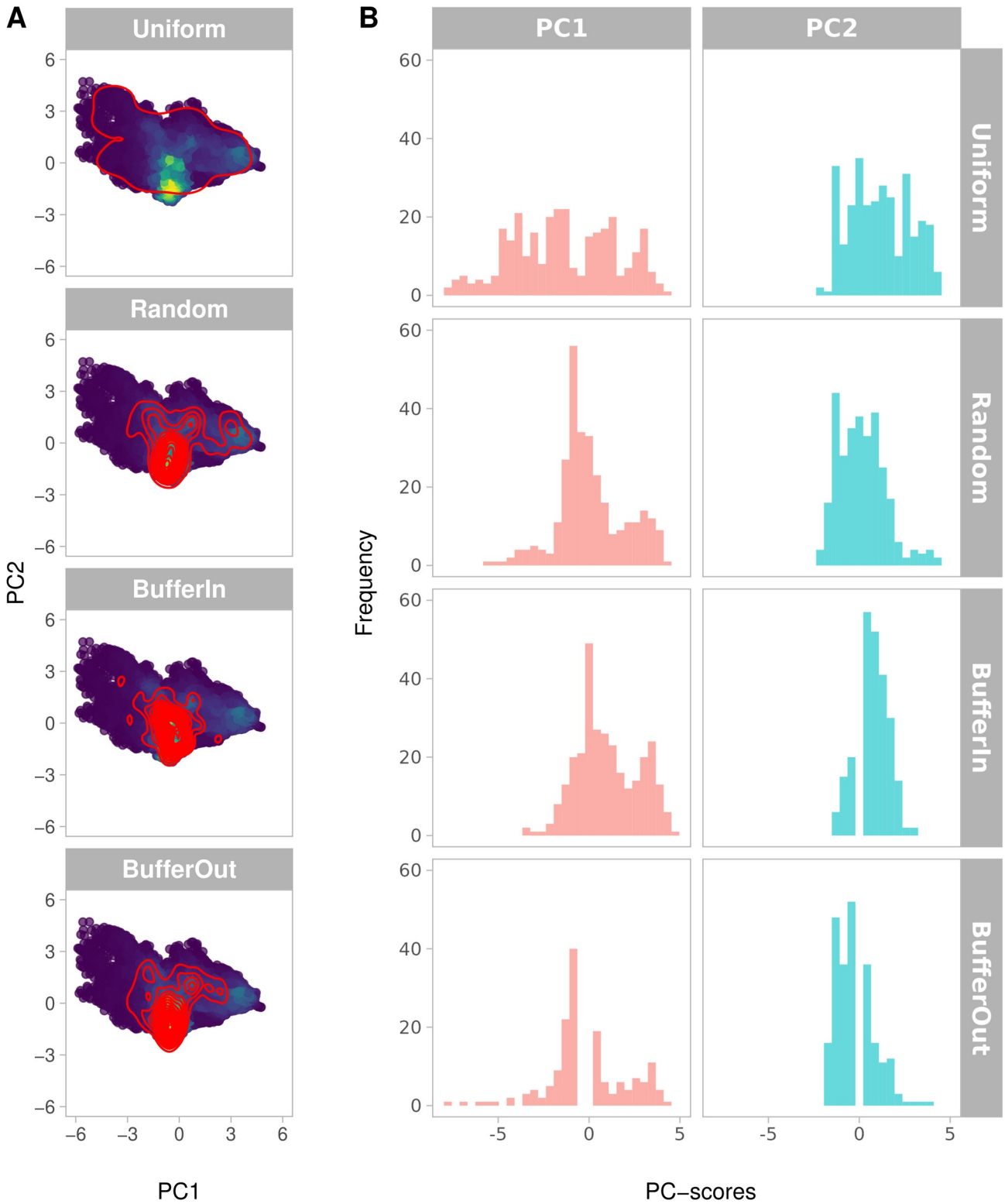
364

## 365  3.2 Effect of sample location bias and class overlap

366  The bivariate density plots of the pseudo-absences sampled within the environmental and
367  geographical space highlighted that the uniform approach had the widest and most
368  homogeneous coverage of environmental conditions throughout the environmental space
369  (Fig. 4, see Figure S2.1  for a detailed overview of the density of pseudo-absences sampled
370  by the uniform approach). In contrast, the random, buffer-in and buffer-out strategies
371  appeared to be prone to sample location bias, with peaks of high density of pseudo-
372  absences occurring in specific areas of the environmental space, i.e. those associated with
373  the most frequent habitat conditions encountered within the geographical space.

374  Regarding class overlap, we detected a significant difference in the overlap between ranges
375  occupied by presence and pseudo-absence points within the environmental space (one-
376  way ANOVA $F(3,36) = 39$, p-value $< 0.001$). Specifically, the uniform approach exhibited the
377  lowest overlap in comparison to the other sampling strategies (Fig. 5).

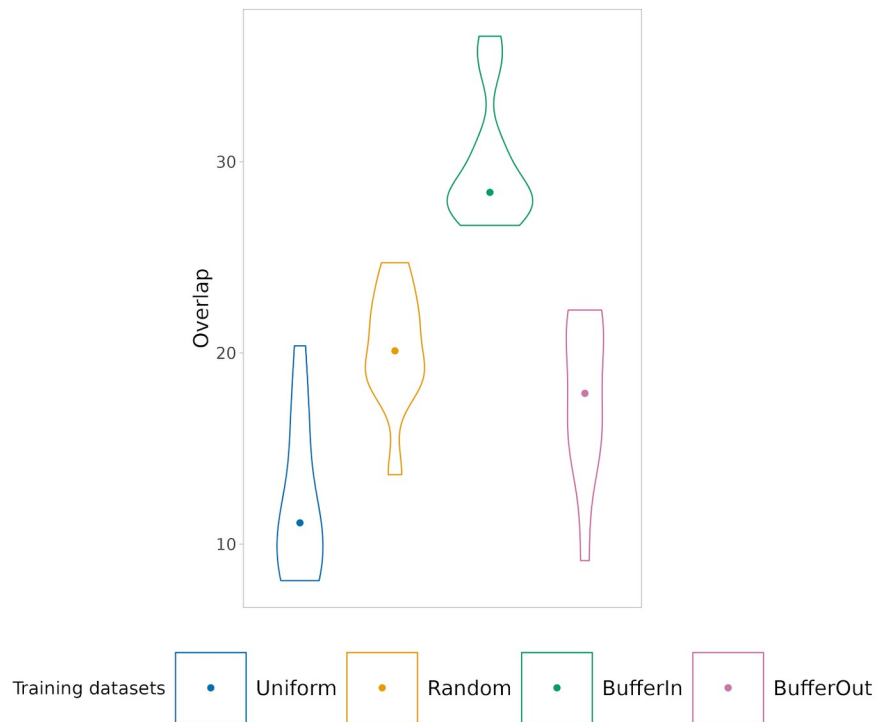USE: a novel approach to uniformly sampling the environmental space

**Figure 4**: A) Bivariate plots showing the environmental space generated by a PCA carried out on 5

bioclimatic variables. Red lines represent the density of pseudo-absences within the environmental

380     space for an individual virtual species. A more detailed representation of the density of pseudo-

381     absences sampled by the uniform approach is available in Figure S2.1. B) Histograms showing the

382     frequency distribution of the first two PCs.



383

384     **Figure 5:** Violin plots showing the overlap in the environmental space between species

385     presences and pseudo-absences. Colours represent samples of pseudo-absences

386     generated using four different strategies, dots represent median values of overlap across 50

387     VS.

388

## 389  4 Discussion

390   In this study, we proposed the uniform approach as an alternative strategy to sample

391   pseudo-absences within the environmental space. In contrast to existing techniques, our

392   approach systematically samples pseudo-absences from portions of the environmental

393   space excluding the conditions that are likely to be suitable for the species to establish. As a

394   result, the uniform approach reduces the chance of including false-absences in the training

395   dataset. From a more theoretical perspective, data collected after the application of the

396   kernel-based filter are much closer to the concept of pseudo-absences than those obtained

397   through traditional, geographical sampling approaches. Our findings showed that the

398   uniform approach represents a valid strategy for gathering pseudo-absences, as it performs

399   equally or better than sampling strategies implemented within the geographical space in

400   terms of model out-of-sample predictive accuracy. Also, the uniform sampling significantly

401   reduces sample location bias and class overlap, which is critical to obtain ecologically

402   meaningful pseudo-absences. Importantly, the uniform approach is flexible, as it lets the

403   user free to set parameters (e.g., kernel bandwidth, sample prevalence, sampling grid

404   resolution) that control how pseudo-absences are sampled within the environmental space.

405   This is particularly valuable, as it makes this approach adaptable for modelling species with

406   different ecological properties (e.g., generalist vs specialist species). By generating

407   informative pseudo-absences, the uniform approach allows satisfying one of the most

408   critical assumptions underpinning habitat suitability modelling: the need of adequate species

409   attributes to model the species-environment relationship (Guisan et al., 2017).

## 410 4.1 Effect of the sampling approaches on models' predictive performances

411 Results of the VS simulations showed that the uniform approach performed well in terms of
412 out-of-sample prediction regardless of the modelling technique, metric of predictive
413 performance, and sample prevalence. HSMs calibrated on pseudo-absences sampled with
414 the uniform approach consistently showed high predictive performance, especially for the
415 accuracy metrics related to the capacity of the model to correctly predict presences (i.e.,
416 sensitivity and CBI). Concerning the metrics associated with the models' ability to predict
417 absences (e.g., specificity), the uniform sampling showed values comparable to the other
418 strategies, except for the buffer-in approach, which always scored the lowest values. This
419 clearly suggests that the uniform approach reduces omission error without necessarily
420 increasing commission error. This is coherent with Fei and Yu (2016), who reported an
421 increase in model predictive performance when pseudo-absences were systematically
422 collected within the environmental space.

423 In this sense, results for the CBI, which is currently the go-to accuracy metric for validating
424 HSMs fitted on pseudo-absences or background points, were particularly encouraging: the
425 uniform approach scored, together with the buffer-out approach, the highest values across
426 all modelling techniques. The high predictive performance associated with the uniform
427 approach can be attributed to how it operates: the systematic sampling of the environmental
428 space and the kernel-based filter. In particular, the specular trends of the predictive
429 accuracy metrics (Fig. 3) and the environmental overlap among pseudo-absences collected
430 through the different sampling approaches (Fig. 5) highlight the importance of the kernel-
431 based filter to favour the discrimination between the environmental features associated with
432 presences and pseudo-absences.

433 Notwithstanding the positive results obtained in terms of predictive performance, we argue

434 that comparing model predictive accuracy may not be the best choice for evaluating the

435 adequacy of sampling carried out within the environmental rather than the geographical

436 space. Indeed, previous studies showed that these metrics are affected by several factors,

437 including sample prevalence (Guisan et al., 2017; Leroy et al., 2018), sample bias (Dubos

438 et al., 2022) or the spatial extent of the study area (Lobo et al., 2008). Moreover, AUC and

439 TSS tend to score high even in case of poor models calibrated on data exhibiting strong

440 sample location bias (Fourcade et al., 2018, Jiménez-Valverde, 2021). Assessing HSMs

441 predictive performance using a set of different predictive accuracy metrics might help the

442 user to critically evaluate the outputs of the models.

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

445 The uniform approach proved to significantly reduce sample location bias, since pseudo-

446 absences were homogeneously scattered along the two principal component axes of the

447 bivariate density plot (Fig. 3ab, Fig S1.2 in Supplementary Materials). On the contrary, the

448 sampling approaches based on geographical space, which all perform a random sampling

449 of the pseudo-absences, exhibited prominent peaks of density of pseudo-absences in

450 correspondence of the most frequently environmental conditions available within the

451 geographical space. As a consequence, the random, buffer-in and buffer-out approaches

452 are likely to provide sub-optimal pseudo-absences for modelling the species-environment

453 relationship, potentially resulting in the estimation of  truncated species response curves

454 (Thuiller et al. 2004; Austin 2007). This aspect is increasingly relevant as environmental

455 conditions are more heterogeneously distributed across the geographic space. Therefore,

456 HSMs calibrated on training datasets adequately representing environmental variability

457  rather than wide geographical coverage represent a crucial step to better capture and

458  discriminate the niche variability of a species (Tessarolo et al., 2014, 2021; Varela et al.,

459  2014; Perret and Sax 2022).

460  The uniform approach proved to also significantly reduce class overlap. The `thres`

461  argument in the `paSampling` function controls the portion of the environmental space

462  associated with the species presence, thus inherently limiting the class overlap issue by

463  excluding environmental conditions more favourable to the study species (see Fig. 1c, 5

464  and Fig. S2.2). This results in a set of pseudo-absences theoretically much closer to the

465  species' true absences. Given that presence points are unevenly distributed within the

466  environmental space, different kernel thresholds might also be used to handle pseudo-

467  absences sampling under particular scenarios. As an example, in case of source-sink

468  dynamics, setting more conservative thresholds for the kernel functions may allow excluding

469  pseudo-absences from environmentally suitable areas, while not excluding areas where a

470  sink population is present due to accidental or mass dispersal close to a source population.

471  ## 4.4 Limitations and usage notes

472  *4.4.1 Limitations*

473  A first limitation of the uniform approach is that its effectiveness depends on the amount

474  (sample size) and quality (e.g., geographically unbiased data *sensu* Fourcade 2014) of

475  presence data. Indeed, if few presence data are available and/or are geographically biased,

476  the kernel-based filter might not accurately delimit the area associated with suitable

477  conditions for the species. As a consequence, the discrimination between suitable and not

478  suitable conditions within the environmental space might be sub-optimal.

479      A second limitation is that, although the uniform approach proved to be robust to

480 varying sample prevalence, its effectiveness might diminish if a very large number of

481 pseudo-absences is sampled (e.g., in case of low sample prevalence) (Fig. S3.1-3.2). Since

482 the uniform approach samples a user-defined number of pseudo-absences within a grid

483 overlaid to the bi-dimensional environmental space, if the number of pseudo-absences

484 grows indefinitely, the advantage of the systematic sampling decreases. Indeed,

485 oversampling the environmental space would generate datasets suffering from sample

486 location bias as much as those based on the random sampling of the geographic space.

487 Finally, from a more practical perspective, the uniform approach can currently only operate

488 across 2-dimensional environmental spaces, but 3-dimensional spaces might be supported

489 in the future.

490 *4.4.2 Usage notes*

491 We here used the uniform approach to sample climatic spaces, although we stress the

492 importance of not only using bioclimatic variables, but also information on soil, land-use as

493 well as other relevant variables when modelling species distributions. Also, we invite

494 potential users of the uniform approach to always check that the first two axes of the

495 principal component analysis used to generate the environmental space explains a large

496 portion of the variance in the data (e.g., ≥ 70%). Equally important is the choice of the

497 boundaries of the geographical extent for which the 2-dimensional space has to be

498 generated. Indeed, to avoid the "there are no elephants in the Antarctic" paradox (Lobo et

499 al., 2010), the spatial extent of the study area should be delineated so that it excludes

500 geographical locations and, in turn, environmental conditions where it is not possible to find

501 the target species due to ecophysiological limitations (e.g., collecting pseudo-absences

502 from Mediterranean coastal dunes when modelling the distribution of an alpine plant

503 species). In short, the uniform approach can provide exhaustive information on where the

504 species is likely to not occur, but it remains a responsibility of the end user to carefully verify

28

USE: a novel approach to uniformly sampling the environmental space

505    if such information is ecologically meaningful.

## 506 5 Conclusion

507 In this study, we evaluated the predictive performance of four sampling strategies, of which one

508 implemented within the environmental space, to collect pseudo-absences for HSMs applications.

509 Also, we compared the sampling approaches in terms of their vulnerability to sample location bias

510 and class overlap. The sampling strategy which we proposed, the uniform approach, proved to (i)

511 have good predictive performances, and (ii) to reduce sample location bias and class overlap. The

512 uniform approach is openly available to users at https://github.com/danddr/USE.

## 513  6 Declaration

514  • Ethics approval and consent to participate: Not applicable.

515  • Competing interests: No conflict of interest has been declared by the authors.

516  • Funding: DDR is supported by a FRS-FNRS Belgian grant, ET is supported by an

517  Estonian Research Council grant (MOBJD1030), MB acknowledges funding from the

518  European Union's Horizon Europe research and innovation programme under the Marie

519  Skłodowska-Curie grant agreement No 101066324.

520  • Authors' contribution: MB conceived the idea of the Uniform approach and wrote the

521  related R functions, while ET and DDR integrated the kernel density-based estimation of

522  presences and the prevalence-related settings. DDR, ET and MB performed the

523  simulations, analysed the data and assembled the USE R package. JL, JJL, SOV, and

524  DR critically commented on the results of the analyses and their interpretation; DDR, ET

525  and MB led the writing of the manuscript and produced a first draft, which was further

526  improved by all other authors.

527  • Acknowledgments: The authors are grateful to Joaquin Hortal, who provided

528  constructive feedback and commented on a previous version of this manuscript.

529  Simulations were carried out using the facilities of the High-Performance Computing

530  Center of the University of Tartu.

## 7 Code and Data availability

531

532 The scripts for replicating the analyses presented in this paper are available at

533 https://github.com/danddr/USE_paper, as well as all the raw outputs of the simulations and

534 statistical analysis, which are available as an .RDS file.

535 We provide a tutorial to explain how to apply the uniform approach to real case studies,

536 using the European beech, *Fagus sylvatica* L. as a target species in S4.  The R script of the

537 tutorial is available at https://github.com/danddr/USE_paper.

# References

Acevedo, P., Jiménez-Valverde, A., Lobo, J. M., and Real, R. (2012). Delimiting the geographical background in species distribution modelling. *Journal of biogeography*, 39(8):1383–1390.

Albert, C. H., Yoccoz, N. G., Edwards Jr, T. C., Graham, C. H., Zimmermann, N. E., and Thuiller, W. (2010). Sampling in ecology and evolution – bridging the gap between theory and practice. *Ecography*, 33(6):1028–1037.

Austin, M. (2007). Species distribution models and ecological theory: a critical assessment and some possible new approaches. *Ecological modelling*, *200*(1-2), 1-19.

Barbet-Massin, M., Jiguet, F., Albert, C. H., and Thuiller, W. (2012). Selecting pseudo absences for species distribution models: how, where and how many? *Methods in ecology and evolution*, 3(2):327–338.

Baker, D. J., Maclean, I. M. D., Goodall, M., & Gaston, K. J. (2022). Correlations between spatial sampling biases and environmental niches affect species distribution models. Global Ecology and Biogeography, 00, 1– 13.

Bazzichetto, M., Massol, F., Carboni, M., Lenoir, J., Lembrechts, J. J., Joly, R., & Renault, D. (2021). Once upon a time in the far south: Influence of local drivers and functional traits on plant invasion in the harsh sub-Antarctic islands. *Journal of Vegetation Science*, *32*(4), e13057.

Beck, J., Böller, M., Erhardt, A., and Schwanghart, W. (2014). Spatial bias in the GBIF database and its effect on modeling species' geographic distributions. *Ecological Informatics*, 19:10–15.

Bedia, J., Herrera, S., and Gutiérrez, J. M. (2013). Dangers of using global bioclimatic datasets for ecological niche modeling. limitations for future climate projections. *Global and Planetary Change*, 107:1–12.

Blonder B, Morrow wcfCB, Harris DJ, Brown S, Butruille G, Laini A, Chen D (2022). _hypervolume: High Dimensional Geometry, Set Operations, Projection, and Inference Using Kernel Density Estimation, Support Vector Machines, and Convex Hulls_. R package version 3.0.4, <https://CRAN.R-project.org/package=hypervolume>.

Booth, T. H., Nix, H. A., Busby, J. R., and Hutchinson, M. F. (2014). Bioclim: the first species distribution modelling package, its early applications and relevance to most current maxent studies. *Diversity and Distributions*, 20(1):1–9.

Da Re, D., Tordoni, E., De Pascalis, F., Negrín-Pérez, Z., Fernández-Palacios, J. M., Arévalo, J. R., ... & Bacaro, G. (2020). Invasive fountain grass (*Pennisetum setaceum* (Forssk.) Chiov.) increases its potential area of distribution in Tenerife island under future climatic scenarios. *Plant Ecology*, *221*(10), 867-882.

Dinno, A. (2017). *dunn.test: Dunn's Test of Multiple Comparisons Using Rank Sums*. R package version 1.3.5, https://CRAN.R-project.org/package=dunn.test.

576 Dubos, N., Préau, C., Lenormand, M., Papuga, G., Monsarrat, S., Denelle, P., ... & Luque,
577   S. (2022). Assessing the effect of sample bias correction in species distribution models.
578   Ecological Indicators, 145, 109487.

579 Duffy, G. A., Coetzee, B. W., Latombe, G., Akerman, A. H., McGeoch, M. A., & Chown, S.
580   L. (2017). Barriers to globally invasive species are weakening across the
581   Antarctic. *Diversity and Distributions*, *23*(9), 982-996.

582 Duong, T. (2021). *ks: Kernel Smoothing*. R package version 1.13.3.

583 Fei, S. and Yu, F. (2016). Quality of presence data determines species distribution model
584 performance: a novel index to evaluate data quality. *Landscape Ecology*, 31(1):31–42.

585 Fick, S. E. and Hijmans, R. J. (2017). Worldclim 2: new 1-km spatial resolution climate
586 surfaces for global land areas. *International journal of climatology*, 37(12):4302–4315.

587 Fourcade, Y. (2021). Fine-tuning niche models matters in invasion ecology. A lesson from
588 the land planarian Obama nungara. *Ecological Modelling*, *457*, 109686.

589 Fourcade, Y., Besnard, A. G., and Secondi, J. (2018). Paintings predict the distribution of
590   species, or the challenge of selecting environmental predictors and evaluation statistics.
591   *Global Ecology and Biogeography*, 27(2):245–256.

592 Fourcade, Y., Engler, J. O., Rödder, D., & Secondi, J. (2014). Mapping species distributions
593   with MAXENT using a geographically biased sample of presence data: a performance
594   assessment of methods for correcting sampling bias. *PloS one*, 9(5), e97122.

595 Grimmett, L., Whitsed, R., & Horta, A. (2020). Presence-only species distribution models
596   are sensitive to sample prevalence: Evaluating models using spatial prediction stability
597   and accuracy metrics. *Ecological Modelling*, *431*, 109194.

598 Guisan, A., Thuiller, W., and Zimmermann, N. E. (2017). *Habitat suitability and distribution
599   models: with applications in R*. Cambridge University Press.

600 Hallgren, W., Santana, F., Low-Choy, S., Zhao, Y., and Mackey, B. (2019). Species
601   distribution models can be highly sensitive to algorithm configuration. *Ecological
602   Modelling*, 408:108719.

603 Hattab, T., Garzón-López, C. X., Ewald, M., Skowronek, S., Aerts, R., Horen, H., Brasseur,
604   B., Gallet-Moron, E., Spicher, F., Decocq, G., et al. (2017). A unified framework to model
605   the potential and realized distributions of invasive species within the invaded range.
606   *Diversity and Distributions*, 23(7):806–819.

607 Hortal, J., Jiménez-Valverde, A., Gómez, J. F., Lobo, J. M., and Baselga, A. (2008).
608   Historical bias in biodiversity inventories affects the observed environmental niche of the
609   species. *Oikos*, 117(6):847–858.

610 Iturbide, M., Bedia, J., Herrera, S., del Hierro, O., Pinto, M., and Gutiérrez, J. M. (2015). A
611   framework for species distribution modelling with improved pseudo-absence generation.
612   *Ecological Modelling*, 312:166–174.

613 Jackson, S. T. and Overpeck, J. T. (2000). Responses of plant populations and
614   communities to environmental changes of the late quaternary. *Paleobiology*, 26(S4):194–

615    220.

616    Jarvie, S., & Svenning, J. C. (2018). Using species distribution modelling to determine
617        opportunities for trophic rewilding under future scenarios of climate change. *Philosophical*
618        *Transactions of the Royal Society B: Biological Sciences*, *373*(1761), 20170446.

619    Jiménez-Valverde, A. (2021). Prevalence affects the evaluation of discrimination capacity in
620        presence-absence species distribution models. *Biodiversity and Conservation*, *30*(5),
621        1331-1340.

622    Jiménez-Valverde, A., Lobo, J. M., & Hortal, J. (2008). Not as good as they seem: the
623        importance of concepts in species distribution modelling. *Diversity and*
624        *distributions*, *14*(6), 885-890.

625    Jiménez-Valverde, A., Acevedo, P., Barbosa, A. M., Lobo, J. M., and Real, R. (2013).
626        Discrimination capacity in species distribution models depends on the representativeness
627        of the environmental domain. *Global Ecology and Biogeography*, 22(4):508–516.

628    Leroy, B., Delsol, R., Hugueny, B., Meynard, C. N., Barhoumi, C., Barbet-Massin, M., and
629        Bellard, C. (2018). Without quality presence–absence data, discrimination metrics such
630        as tss can be misleading measures of model performance. *Journal of Biogeography*,
631        45(9):1994–2002.

632    Leroy, B., Meynard, C. N., Bellard, C., and Courchamp, F. (2016). virtualspecies, an r
633        package to generate virtual species distributions. *Ecography*, 39(6):599–607.

634    Lobo, J. M., Jiménez-Valverde, A., and Hortal, J. (2010). The uncertain nature of absences
635        and their importance in species distribution modelling. *Ecography*, 33(1):103–114.

636    Lobo, J. M., Jiménez-Valverde, A., and Real, R. (2008). Auc: a misleading measure of the
637        performance of predictive distribution models. *Global ecology and Biogeography*,
638        17(2):145– 151.

639    Magri, D., Vendramin, G. G., Comps, B., Dupanloup, I., Geburek, T., Gömöry, D., ... & De
640        Beaulieu, J. L. (2006). A new scenario for the Quaternary history of European beech
641        populations: palaeobotanical evidence and genetic consequences. *New phytologist*,
642        171(1), 199-221.

643    Mammola, S. and Cardoso, P. (2020). Functional diversity metrics using kernel density n-
644        dimensional hypervolumes. *Methods in Ecology and Evolution*, 11(8):986–995.

645    Meynard, C. N., Leroy, B., and Kaplan, D. M. (2019). Testing methods in species
646        distribution modelling using virtual species: what have we learnt and what are we
647        missing? *Ecography*, 42(12):2021–2036.

648    Naimi, B. and Araújo, M. B. (2016). sdm: a reproducible and extensible r platform for
649        species distribution modelling. *Ecography*, 39(4):368–375.

650    Newbold, T. (2018). Future effects of climate and land-use change on terrestrial vertebrate
651        community diversity under different scenarios. *Proceedings of the Royal Society B*,
652        285(1881):20180792.

653    Perret, D. L. and Sax, D. F. (2022). Evaluating alternative study designs for optimal

654  sampling of species' climatic niches. *Ecography*.

655  Poli et al. (2022) Coupling fossil records and traditional discrimination metrics to test how
656  genetic information improves species distribution models of the European beech Fagus
657  sylvatica. *European Journal of Forest Research*, 141: 253–265

658  Phillips, S. J., Anderson, R. P., Dudík, M., Schapire, R. E., and Blair, M. E. (2017). Opening
659  the black box: An open-source release of maxent. *Ecography*, 40(7):887–893.
660  Phillips, S. J., Dudík, M., Elith, J., Graham, C. H., Lehmann, A., Leathwick, J., and Ferrier,
661  S. (2009). Sample selection bias and presence-only distribution models: implications for
662  background and pseudo-absence data. *Ecological applications*, 19(1):181–197.

663  R Core Team (2021). *R: A Language and Environment for Statistical Computing*. R
664  Foundation for Statistical Computing, Vienna, Austria.

665  Ronquillo, C., Alves-Martins, F., Mazimpaka, V., Sobral-Souza, T., Vilela-Silva, B., Medina,
666  N. G., and Hortal, J. (2020). Assessing spatial and temporal biases and gaps in the
667  publicly available distributional information of iberian mosses. *Biodiversity Data Journal*,
668  8.

669  Santini, L., Benítez-López, A., Maiorano, L., Čengić, M., and Huijbregts, M. A. (2021).
670  Assessing the reliability of species distribution projections in climate change research.
671  *Diversity and Distributions*, 27(6):1035–1050.

672  Sillero, N. and Barbosa, A. M. (2020). Common mistakes in ecological niche models.
673  *International Journal of Geographical Information Science*, pages 1–14.

674  Støa, B., Halvorsen, R., Stokland, J. N., and Gusarov, V. I. (2019). How much is enough?
675  influence of number of presence observations on the performance of species distribution
676  models. *Sommerfeltia*, 39(1):1–28.

677  Svenning, J.-C. and Skov, F. (2004). Limited filling of the potential range in European tree
678  species. *Ecology Letters*, 7(7):565–573.

679  Tessarolo, G., Lobo, J. M., Rangel, T. F., and Hortal, J. (2021). High uncertainty in the
680  effects of data characteristics on the performance of species distribution models.
681  *Ecological Indicators*, 121:107147.

682  Tessarolo, G., Rangel, T. F., Araújo, M. B., and Hortal, J. (2014). Uncertainty associated
683  with survey design in species distribution models. *Diversity and Distributions*,
684  20(11):1258–1269.

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

688  Valavi, R., Elith, J., Lahoz-Monfort, J. J., and Guillera-Arroita, G. (2021). Modelling species
689  presence-only data with random forests. *Ecography*, 44(12):1731–1742.

690  Varela, S., Anderson, R. P., García-Valdés, R., and Fernández-González, F. (2014).
691  Environmental filters reduce the effects of sampling bias and improve predictions of
692  ecological niche models. *Ecography*, 37(11):1084–1091.

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

695   Wasof et al. (2015) Disjunct populations of European vascular plant species keep the same
696      climatic niches. *Global Ecology and Biogeography*, 24: 1401-1412

697   Wickham H (2016). ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York. ISBN
698   978-3-319-24277-4, https://ggplot2.tidyverse.org.

# Supplementary Material 1

**Tab. S1**: Post-hoc multiple comparisons with Dunn's rank sum test (α = 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 | χ2 | P.val |
|-------|--------|-------------|-----|-------|
| BRT | AUC | Buffer IN - Uniform | 146.28 | p<0.001 |
| BRT | AUC | Buffer OUT - Uniform | 146.28 | 0.176 |
| BRT | AUC | Random - Uniform | 146.28 | p<0.001 |
| BRT | BoyceI | Buffer IN - Uniform | 131.66 | p<0.001 |
| BRT | BoyceI | Buffer OUT - Uniform | 131.66 | 0.0356 |
| BRT | BoyceI | Random - Uniform | 131.66 | 0.009 |
| BRT | Sensitivity | Buffer IN - Uniform | 104.47 | p<0.001 |
| BRT | Sensitivity | Buffer OUT - Uniform | 104.47 | p<0.001 |
| BRT | Sensitivity | Random - Uniform | 104.47 | p<0.001 |
| BRT | Specificity | Buffer IN - Uniform | 76.62 | p<0.001 |
| BRT | Specificity | Buffer OUT - Uniform | 76.62 | 0.016 |
| BRT | Specificity | Random - Uniform | 76.62 | p<0.001 |
| BRT | TSS | Buffer IN - Uniform | 150.78 | p<0.001 |
| BRT | TSS | Buffer OUT - Uniform | 150.78 | 0.0087 |
| BRT | TSS | Random - Uniform | 150.78 | p<0.001 |
| GAM | AUC | Buffer IN - Uniform | 141.11 | p<0.001 |
| GAM | AUC | Buffer OUT - Uniform | 141.11 | 0.0336 |
| GAM | AUC | Random - Uniform | 141.11 | p<0.001 |
| GAM | BoyceI | Buffer IN - Uniform | 144.02 | p<0.001 |
| GAM | BoyceI | Buffer OUT - Uniform | 144.02 | 0.0044 |
| GAM | BoyceI | Random - Uniform | 144.02 | 0.0033 |
| GAM | Sensitivity | Buffer IN - Uniform | 131.32 | p<0.001 |
| GAM | Sensitivity | Buffer OUT - Uniform | 131.32 | p<0.001 |
| GAM | Sensitivity | Random - Uniform | 131.32 | p<0.001 |

| Model | Metric | Comparisons | χ2 | P.val |
|-------|--------|-------------|-----|-------|
| GAM | Specificity | Buffer IN - Uniform | 128.72 | p<0.001 |
| GAM | Specificity | Buffer OUT - Uniform | 128.72 | 0.1586 |
| GAM | Specificity | Random - Uniform | 128.72 | p<0.001 |
| GAM | TSS | Buffer IN - Uniform | 145.45 | p<0.001 |
| GAM | TSS | Buffer OUT - Uniform | 145.45 | 0.0028 |
| GAM | TSS | Random - Uniform | 145.45 | p<0.001 |
| GLM | AUC | Buffer IN - Uniform | 132.53 | p<0.001 |
| GLM | AUC | Buffer OUT - Uniform | 132.53 | 0.0822 |
| GLM | AUC | Random - Uniform | 132.53 | 0.003 |
| GLM | Boycel | Buffer IN - Uniform | 175.57 | p<0.001 |
| GLM | Boycel | Buffer OUT - Uniform | 175.57 | p<0.001 |
| GLM | Boycel | Random - Uniform | 175.57 | p<0.001 |
| GLM | Sensitivity | Buffer IN - Uniform | 128.02 | p<0.001 |
| GLM | Sensitivity | Buffer OUT - Uniform | 128.02 | p<0.001 |
| GLM | Sensitivity | Random - Uniform | 128.02 | p<0.001 |
| GLM | Specificity | Buffer IN - Uniform | 98.02 | p<0.001 |
| GLM | Specificity | Buffer OUT - Uniform | 98.02 | p<0.001 |
| GLM | Specificity | Random - Uniform | 98.02 | 0.1366 |
| GLM | TSS | Buffer IN - Uniform | 141.06 | p<0.001 |
| GLM | TSS | Buffer OUT - Uniform | 141.06 | 0.0333 |
| GLM | TSS | Random - Uniform | 141.06 | p<0.001 |
| Maxent | AUC | Buffer IN - Uniform | 151.46 | p<0.001 |
| Maxent | AUC | Buffer OUT - Uniform | 151.46 | 0.0099 |
| Maxent | AUC | Random - Uniform | 151.46 | p<0.001 |
| Maxent | Boycel | Buffer IN - Uniform | 178.36 | p<0.001 |
| Maxent | Boycel | Buffer OUT - Uniform | 178.36 | p<0.001 |
| Maxent | Boycel | Random - Uniform | 178.36 | p<0.001 |
| Maxent | Sensitivity | Buffer IN - Uniform | 64.45 | p<0.001 |
| Maxent | Sensitivity | Buffer OUT - Uniform | 64.45 | 0.0677 |

| Model | Metric | Comparisons | χ2 | P.val |
|---|---|---|---|---|
| Maxent | Sensitivity | Random - Uniform | 64.45 | 0.0099 |
| Maxent | Specificity | Buffer IN - Uniform | 66.81 | p<0.001 |
| Maxent | Specificity | Buffer OUT - Uniform | 66.81 | 0.046 |
| Maxent | Specificity | Random - Uniform | 66.81 | 0.0035 |
| Maxent | TSS | Buffer IN - Uniform | 151.49 | p<0.001 |
| Maxent | TSS | Buffer OUT - Uniform | 151.49 | 0.0098 |
| Maxent | TSS | Random - Uniform | 151.49 | p<0.001 |
| RF | AUC | Buffer IN - Uniform | 147.3 | p<0.001 |
| RF | AUC | Buffer OUT - Uniform | 147.3 | 0.0747 |
| RF | AUC | Random - Uniform | 147.3 | p<0.001 |
| RF | Boycel | Buffer IN - Uniform | 166.26 | p<0.001 |
| RF | Boycel | Buffer OUT - Uniform | 166.26 | 0.1462 |
| RF | Boycel | Random - Uniform | 166.26 | p<0.001 |
| RF | Sensitivity | Buffer IN - Uniform | 89.75 | p<0.001 |
| RF | Sensitivity | Buffer OUT - Uniform | 89.75 | p<0.001 |
| RF | Sensitivity | Random - Uniform | 89.75 | 0.1444 |
| RF | Specificity | Buffer IN - Uniform | 108.22 | p<0.001 |
| RF | Specificity | Buffer OUT - Uniform | 108.22 | p<0.001 |
| RF | Specificity | Random - Uniform | 108.22 | p<0.001 |
| RF | TSS | Buffer IN - Uniform | 147.11 | p<0.001 |
| RF | TSS | Buffer OUT - Uniform | 147.11 | 0.079 |
| RF | TSS | Random - Uniform | 147.11 | p<0.001 |

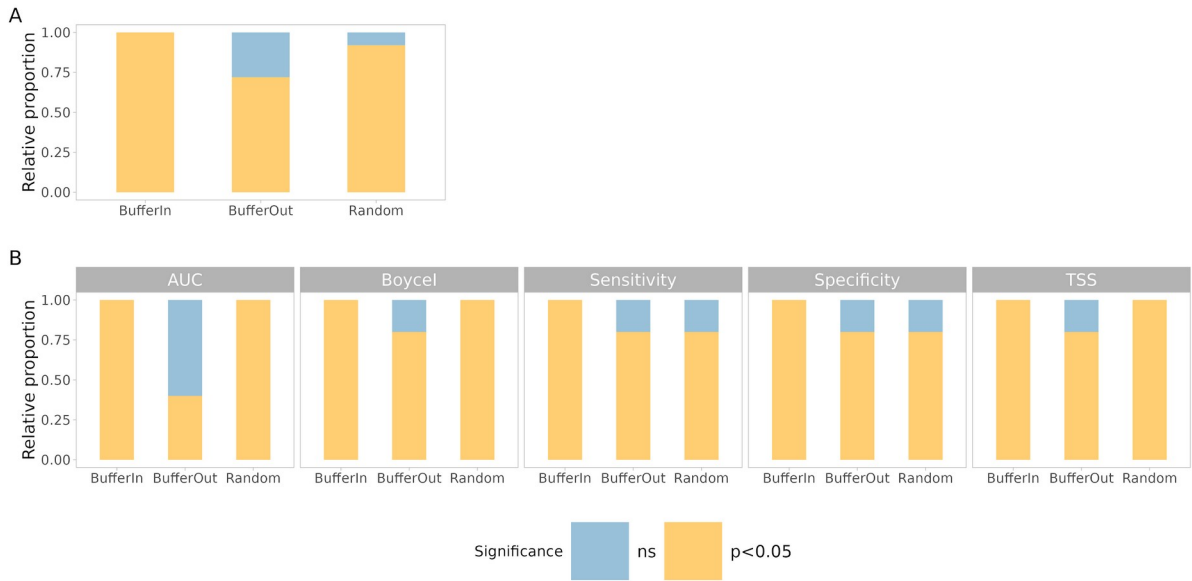**Figure S1.1**: Post-hoc multiple comparisons with Dunn's rank sum test (α = 0.05; omnibus test was always significant with P < 0.05, data not shown). All the comparisons were performed comparing the Uniform dataset with the other different sampling strategies: A) relative proportion of the significant comparisons aggregated by sampling strategy; B) relative proportion of the significant comparisons aggregated by sampling strategy and metric.
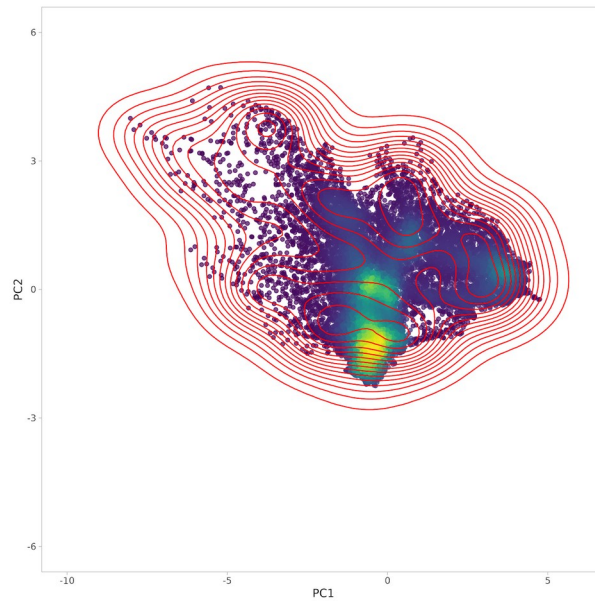
# Supplementary Material 2



**Figure S2.1**: Bivariate plots showing the environmental space generated by a PCA carried out on 5 bioclimatic variables. Red lines represent the density of pseudo-absences within the environmental space for an individual virtual species and for the uniform approach only.
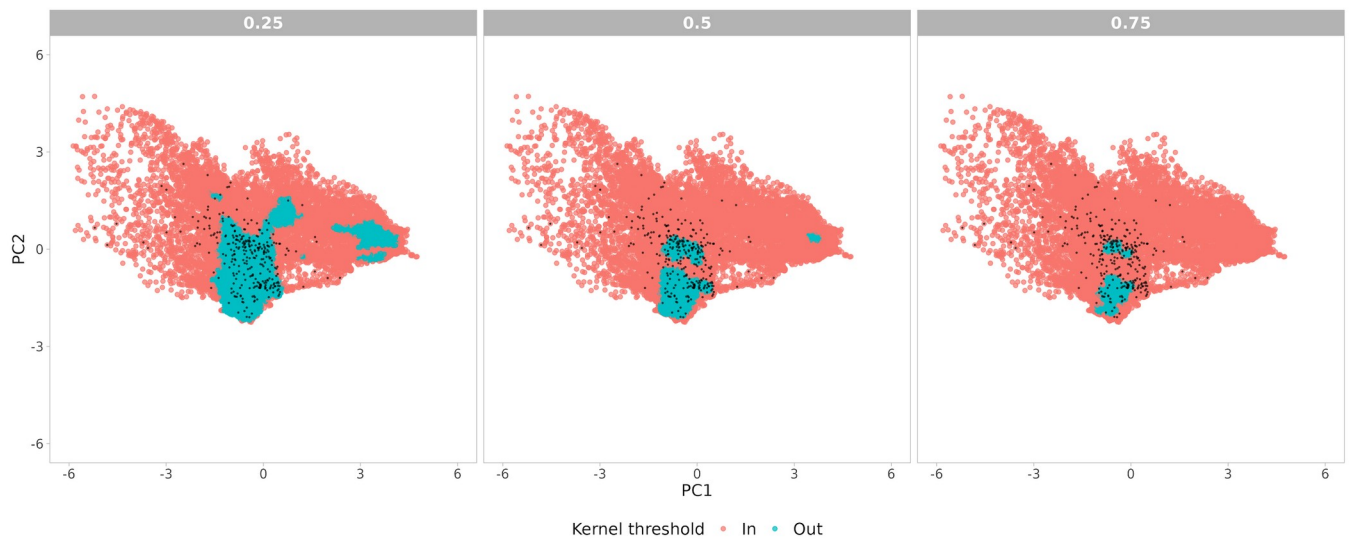


**Figure S2.2**: Effect of the kernel threshold in the inclusion/exclusion of the environmental space to sample. Black dots are the real VS occurrences plotted in the environmental space.

# Supplementary Material 3

## Sensitivity analyses on the sample prevalence

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).
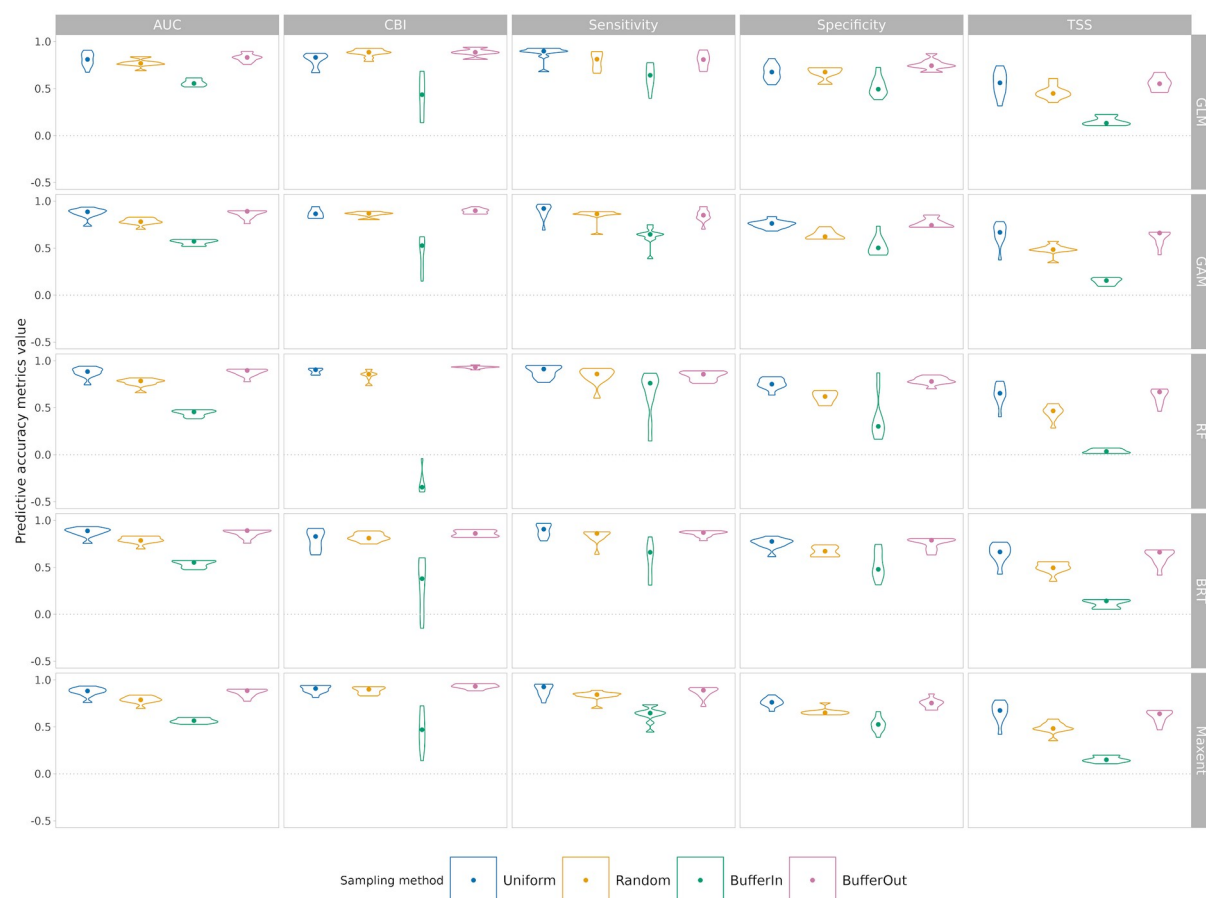


**Figure S3.1**: Violin plots reporting the distribution of the values of the metrics of predictive performance for the HSMs of the 50 VS  modelled as a function of 5, randomly selected bioclimatic predictors, and setting sample prevalence equal to 0.5. Dots represent median values of the metrics of predictive accuracy. Columns indicate the different performance metrics, while rows the modelling techniques used to compute HSMs.
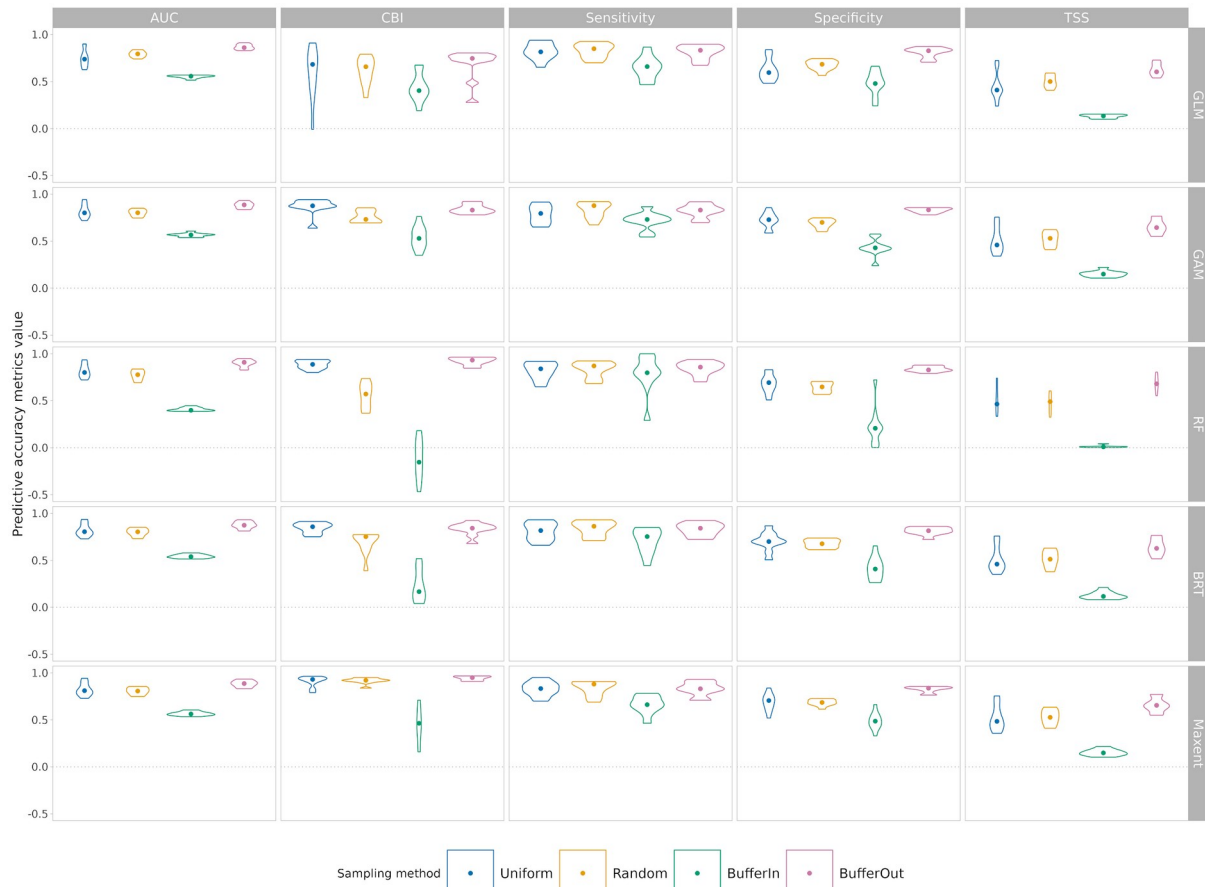
**Figure S3.2**: Violin plots reporting the distribution of the values of the metrics of predictive performance for the HSMs of the 50 VS  modelled as a function of 5, randomly selected bioclimatic predictors, and setting sample prevalence equal to 0.1. Dots represent median values of the metrics of predictive accuracy. Columns indicate the different performance metrics, while rows the modelling techniques used to compute HSMs.

# Supplementary Material 4: case study on the realized distribution of *Fagus sylvatica* in Western Europe

## Methods

To illustrate how to apply the uniform approach using 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 only two modelling algorithms (i.e., GLM and RF). Indeed, the case study of *F. sylvatica* is only used as a practical example on how to use the `USE` package, while not providing a further comparison of the predictive performance of HSMs fitted on data collected through different sampling strategies (as already done with the VS approach). We gathered data on the presence-absence of *F. sylvatica* from the open EU-Forest dataset (Mauri et al., 2017), which compiles presence data 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-absence records of tree species exhaustively collected across Europe, and then aggregated to a 1 × 1 km resolution grid. This let us assume with a certain degree of confidence that the EU-Forest dataset provided a geographically unbiased sample of occurrence records for *F. sylvatica*, and absence data represented 'true' absences.

Across our study area, the EU-Forest dataset provided 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 presence (for internal validation) dataset. To this aim, we generated a 2-dimensional environmental space using all 19 bioclimatic variables available from WorldClim. Then, we used the function `uniformSampling` from the `USE` package to uniformly sample occurrence records within the environmental space. Note that this approach is conceptually similar to the spatial-thinning approach 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 we here applied it uniformly and within the environmental space (see Varela et al., 2014;

Hattab et al. 2017). Once presence records (n = 2,747) were uniformly sub-sampled within the environmental space, we randomly divided them into two sets of training (70%) and testing (30%) occurrences to then derive the two respective sets of training and testing pseudo-absences. To this aim, we first used all 12,444 available presence records to recover the core area of *F. sylvatica*'s bioclimatic niche using the function `paSampling` from the `USE` package and then filtered out the background points likely associated with suitable locations for the species (see step 1 in section 2.2.1 of the main text). Once we removed background points likely associated with the core bioclimatic niche of *F. sylvatica*, the obtained sample sizes were: 1,856 and 906 background points for the training and testing (internal validation) dataset, respectively. Finally, we derived a completely independent testing (external validation) dataset using presence and true absence data from sPlotOpen (Sabatini et al., 2021). The sPlotOpen database 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 ($n$ = 366), and also to derive true absence data from those vegetation plots where *F. sylvatica* was not recorded ($n$ = 4039). As done for the EU-Forest dataset, we selected only those vegetation plots data from sPlotOpen included in our study area (Italy, France and Spain) in western Europe.

The realised distribution of *F. sylvatica* was modelled as a function of WorldClim bioclimatic variables (resolution of 2.5 minutes at the Equator). For simplicity, we solely focused on the climatic niche of *Fagus sylvatica*, although we acknowledge that other factors different from climatic drivers may equally contribute in shaping the distribution of this species, especially so at local scales (Mellert et al., 2018). As modelling techniques, we used binary generalised linear models and random forests (ranger function available in ranger R package; Wright and Ziegler, 2017). To avoid multicollinearity issues, we selected a subset of the 19 bioclimatic variables using the findCorrelation function from the caret R package (Kuhn, 2021) (pairwise-correlation threshold: 0.6). The bioclimatic variables finally selected for *F. sylvatica* were: BIO6 (minimum temperature of the coldest month); BIO7 (temperature annual range); and BIO8 (mean temperature of the wettest quarter). We also used the latitudinal position of the presence and pseudo-absence records as an additional predictor to account for the effect of factors affecting *F. sylvatica* that correlates with its latitudinal gradient of

occurrence and were not included in the model, such as its 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 binary GLM, we introduced second order polynomial terms for BIO6, BIO7 and latitude. Statistically non-significant predictors were dropped from the original full model to reach a most parsimonious model. The predictive performance of the fitted models was assessed using TSS and CBI on three different types of data: (i) the 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 testing dataset derived from sPlotOpen. We often assume TSS>0.5 to indicate good predictions, while CBI positive values indicate a model which present predictions are consistent with the distribution of presences in the evaluation dataset, values close to zero mean that the model is not different from a random model, negative values indicate counter predictions, i.e., predicting poor quality areas where presences are more frequent (Hirzel et al. 2006).

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

A full description of the modelling procedure (from the sub-sampling of the presence and background points to the assessment of the model predictive performance) is reported at https://github.com/danddr/USE_paper.

## Results

Both the GLM and the RF for *F. sylvatica* showed high predictive performances, regardless of the dataset used for testing (Table 1). Concerning the GLM, the TSS was always equal to or above 0.40, with the lowest value obtained for the sPlotOpen testing dataset (0.40) and the highest for the EU-Forest dataset (0.6). The lowest CBI was scored for the sPlotOpen dataset (0.88), while the highest for the EU-Forest dataset (0.98).

Similar results were obtained for the RF, for which the lowest TSS was obtained for the sPlotOpen testing dataset (0.48), while comparable values resulted from the EU-Forest dataset and the 5-fold cross validation (0.79 and 0.76, respectively). With respect to the CBI, the highest value was observed for the EU-Forest dataset (average = 0.98), while the lowest was obtained for the sPlotOpen dataset (0.96).

Goodness-of-fit measures seemed to be affected by the modelling technique, with the $R^2$ of the RF being 0.67, and the Tjur's $R^2$ for the GLM being 0.35 (Tab. S4.1).

The pseudo-absences of *F. sylvatica* collected using the uniform approach are uniformly distributed in the environmental space, reducing the sample location bias (Fig. S4.1a) and the class overlap, since in the geographical space they are distributed in areas far from the occurrences (Fig. S4.1b).

**Table S4.1**: Results of the two 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; Boyce I: Boyce's Index; R-sq: Tjur's $R^2$ for the GLM, and $R^2$ for RF.

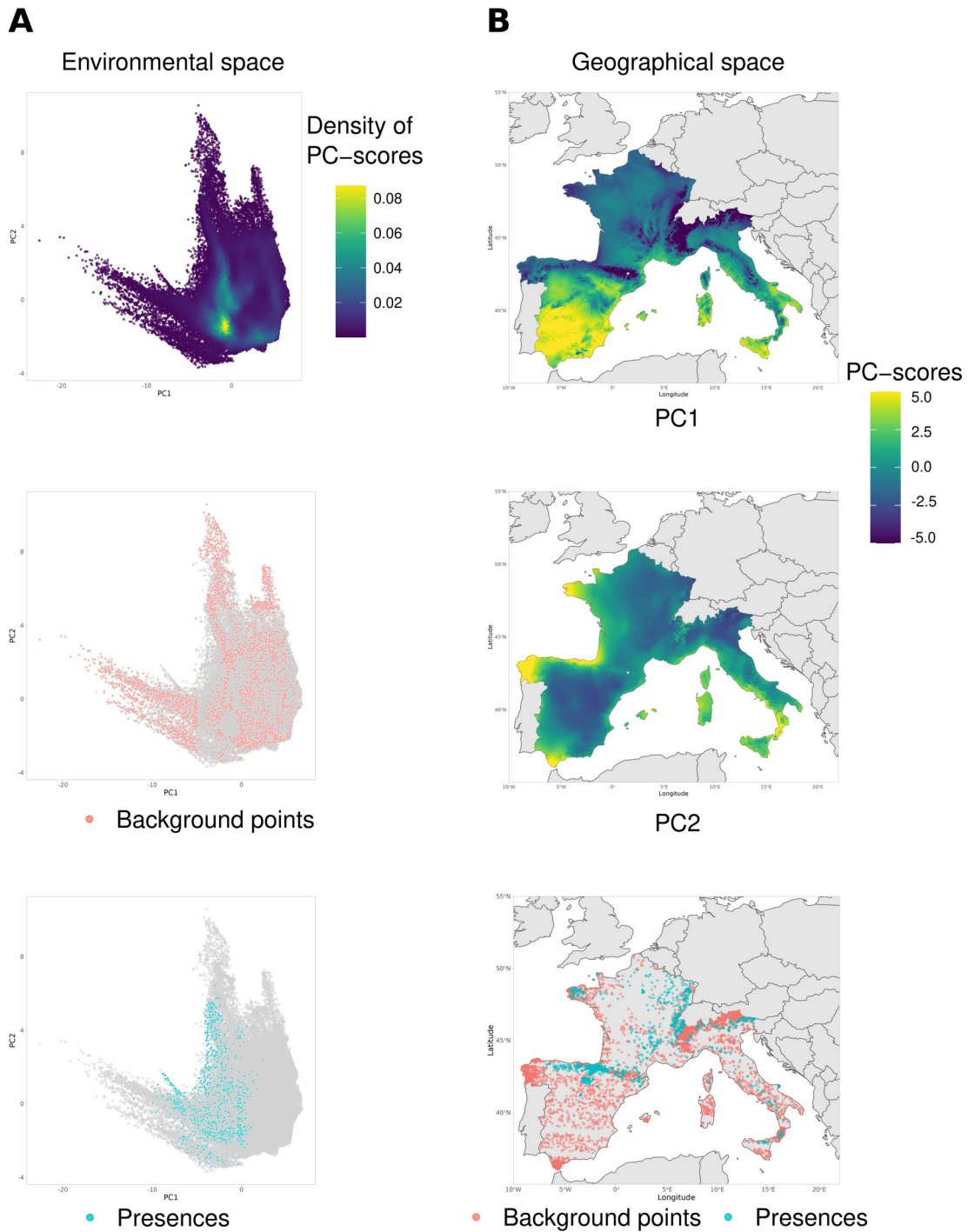| Validation dataset | GLM | | | RF | | |
|---|---|---|---|---|---|---|
| | *TSS* | *CBI* | *Tjur's $R^2$* | *TSS* | *CBI* | *$R^2$* |
| *5-fold CV* | 0.52 | 0.92 | | 0.76 | 0.97 | |
| *EU-Forest* | 0.60 | 0.98 | 0.35 | 0.79 | 0.98 | 0.67 |
| *sPlotOpen* | 0.40 | 0.88 | | 0.48 | 0.96 | |

**Figure S4.1:** (A) the environmental space available for *Fagus sylvatica* in Italy, Spain and France, and the position of presences and pseudo-absences sampled within the environmental space using the Uniform approach; (B) the distribution of PC-scores in the geographical space and the geographical location of presences (blue) and pseudo-absences (red) sampled within the environmental space using the Uniform approach.

# References

Aiello-Lammens, M. E., Boria, R. A., Radosavljevic, A., Vilela, B., and Anderson, R. P. (2015). spthin: an R package for spatial thinning of species occurrence records for use in ecological niche models. *Ecography*, 38(5):541–545.

Hattab, T., Garzón-López, C. X., Ewald, M., Skowronek, S., Aerts, R., Horen, H., Brasseur, B., Gallet-Moron, E., Spicher, F., Decocq, G., et al. (2017). A unified framework to model the potential and realized distributions of invasive species within the invaded range. *Diversity and Distributions*, 23(7):806–819.

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., Chytr`y, 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.

Varela, S., Anderson, R. P., García-Valdés, R., and Fernández-González, F. (2014). Environmental filters reduce the effects of sampling bias and improve predictions of ecological niche models. *Ecography*, 37(11):1084–1091.

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.