

Forest tree species distribution for Europe 2000-2020: mapping potential and realized distributions using spatiotemporal Machine Learning

Carmelo Bonannella (✉ carmelo.bonannella@opengeohub.org)

OpenGeoHub <https://orcid.org/0000-0002-5391-8427>

Tomislav Hengl

OpenGeoHub <https://orcid.org/0000-0002-9921-5129>

Johannes Heisig

University of Münster <https://orcid.org/0000-0003-3586-3001>

Leandro Parente

OpenGeoHub <https://orcid.org/0000-0003-1589-0467>

Marvin N Wright

University of Bremen <https://orcid.org/0000-0002-8542-6291>

Martin Herold

Wageningen University & Research <https://orcid.org/0000-0003-0246-6886>

Sytze de Bruin

Wageningen University & Research <https://orcid.org/0000-0002-6884-2832>

Research Article

Keywords: Species distribution model, Ecological niche, Ensemble modeling, Machine learning, Imbalanced data, Presence-absence, Spatiotemporal modeling, Stacked generalization, High resolution, Tree species

Posted Date: January 12th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1252972/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Forest tree species distribution for Europe**
2 **2000–2020: mapping potential and realized**
3 **distributions using spatiotemporal Machine**
4 **Learning**

5 **Carmelo Bonannella¹, Tomislav Hengl¹, Johannes Heisig², Leandro**
6 **Parente¹, Marvin N Wright^{3, 4}, Martin Herold^{5, 6}, and Sytze de Bruin⁶**

7 ¹**OpenGeoHub, Wageningen, The Netherlands**

8 ²**Institute for Geoinformatics, University of Münster, Münster, Germany**

9 ³**Leibniz Institute for Prevention Research and Epidemiology – BIPS, Bremen, Germany**

10 ⁴**University of Bremen, Bremen, Germany**

11 ⁵**GFZ German Research Centre for Geosciences, Section 1.4 Remote Sensing and**
12 **Geoinformatics, Telegrafenberg, Potsdam, 14473, Germany**

13 ⁶**Laboratory of Geo-Information Science and Remote Sensing, Wageningen University**
14 **& Research, Wageningen, The Netherlands**

15 Corresponding author:

16 Carmelo Bonannella¹

17 Email address: carmelo.bonannella@opengeohub.org

18 **ABSTRACT**

19 Paper describes a data-driven framework based on spatio-temporal ensemble machine learning to produce
20 distribution maps for 16 forest tree species (*Abies alba* Mill., *Castanea sativa* Mill., *Corylus avellana* L.,
21 *Fagus sylvatica* L., *Olea europaea* L., *Picea abies* L. H. Karst., *Pinus halepensis* Mill., *Pinus nigra* J. F.
22 Arnold, *Pinus pinea* L., *Pinus sylvestris* L., *Prunus avium* L., *Quercus cerris* L., *Quercus ilex* L., *Quercus*
23 *robur* L., *Quercus suber* L. and *Salix caprea* L.) at high spatial resolution (30 m). Tree occurrence data
24 for a total of 3 million of points was used to train different Machine Learning (ML) algorithms: random
25 forest, gradient-boosted trees, generalized linear models, k-nearest neighbors, CART and an artificial neural
26 network. A stack of 585 coarse and high resolution covariates representing spectral reflectance (Landsat
27 bands, spectral indices; time-series of seasonal composites), different biophysical conditions (i.e. temperature,
28 precipitation, elevation, lithology) and biotic competition (other species distribution maps) was used as
29 predictors for realized distributions, while potential distribution was modelled with environmental predictors
30 only. Logloss and computing time were used to select the three best algorithms to train an ensemble model
31 based on stacking with a logistic regressor as a meta-learner for each species. High resolution (30 m)
32 probability and model uncertainty maps of realized distribution were produced for each species using a time
33 window of 4 years for a total of 6 distribution maps per species for the studied period, while for potential
34 distributions only one map per species was produced. Results of spatial cross validation show that *Olea*
35 *europaea* and *Quercus suber* achieved the best performances in both potential and realized distribution,
36 while *Pinus sylvestris* and *Salix caprea* achieved the worst. Further analysis shows that fine-resolution models
37 consistently outperformed coarse resolution models (250 m) for realized distribution (average decrease
38 in logloss: +53%). Realized distribution models achieved higher predictive performances than potential
39 distribution ones. Importance of predictor variables differed across species and models, with the green
40 band for summer and the NDWI and NDVI for fall for realized distribution and the diffuse irradiation and
41 precipitation of the driest quarter being the most important and frequent for potential distribution. The
42 ensemble model outperformed or performed as good as the best individual model in all potential species
43 distributions, while for ten species it performed worse than the best individual model in modeling realized
44 distributions. The framework shows how combining continuous and consistent EO time series data with
45 state of the art ML can be used to derive dynamic distribution maps. The produced time-series occurrence
46 predictions can be used to quantify temporal trends and detect potential forest degradation.

47 Submitted to PeerJ on: 7th of January 2022;

48 1 INTRODUCTION

49 Reforestation and restoration are considered key strategies for tackling global warming and enhancing
50 CO₂ sequestration (Lefebvre et al., 2021; Domke et al., 2020; Nave et al., 2019). Under the European
51 Green Deal and the EU biodiversity strategy for 2030, the EU has committed to planting at least 3 billion
52 additional trees in the EU by 2030 (https://ec.europa.eu/environment/3-billion-trees_en).
53 At the same time in Germany, trees in more than 2.5% of the country's total forest area have died because of
54 beetles and drought fueled by a warming climate (Popkin, 2021). Obtaining reliable information on forest
55 tree species distribution in space and time has become paramount for stakeholders and decision-makers to
56 anticipate and minimize climate change impacts (Keenan, 2015).

57 The production of species distribution maps has been dominated by Species Distribution Modeling

58 (SDM) approaches (Franklin, 2010). Maps of species ecological niches are made by associating values
59 of different predictors to known locations of the target species and then used to predict distribution in
60 geographic space where no field data for the target species is available. Correlative SDM (C-SDM) in
61 particular has recently become very popular due to the increase in advanced statistical techniques (i.e.
62 Machine Learning) and geospatial data (satellites, drones etc.) availability (Gobeyn et al., 2019). This
63 is also reflected in an increase of SDM-related publications in the last decade (Booth, 2018), as well as
64 reviews (Booth et al., 2014).

65 From the late 2000s, the focus in SDM has shifted from traditional statistical models to machine-
66 learning (ML) algorithms (Elith et al., 2008). While statistical models start by assuming a particular
67 structural model whose parameters are estimated from the data, ML tries to learn the relationship between
68 the response and the predictors through the observation of dominant patterns (Breiman, 2001b). There
69 are no ecological assumptions explicitly embedded in the algorithms (e.g. Random Forest, Artificial
70 Neural Networks) and most of the time it is difficult to interpret their outputs, which is also referred to
71 as the “black box” problem of ML (Molnar, 2020). However, the exponential increase in computing
72 power (Gorelick et al., 2017), predictor variables availability (Zhu et al., 2019) and ecological “big data”
73 gathered by multiple sources (i.e. sensors, cameras etc.) (Hampton et al., 2013) make ML increasingly
74 suitable tool for SDM. Possibly it is even the only practical means, to gain insights from such an amount
75 of information and to map, monitor or forecast changes on multiple geographical scales (Gobeyn et al.,
76 2019).

77 There is a broad variety of algorithms used for SDM (Gobeyn et al., 2019; Franklin, 2010): among the
78 most used presence-absence we find Generalized Linear Models (GLM) (Nelder and Wedderburn, 1972),
79 Classification and Regression Trees (CART) (Quinlan, 1986), gradient-boosted Trees (GBT) (Friedman,
80 2002), Random Forest (RF) (Breiman, 2001a) and Artificial Neural Networks (ANN) (McCulloch and
81 Pitts, 1943). In literature, many examples of tree species distribution maps created with regression-based
82 and ML SDM can be found: Hill et al. (2017) used GLM to model the distribution of *Acer platanoides* and
83 *Fraxinus excelsior* for Great Britain at 1 km resolution, Marchi and Ducci (2018) compared the predictive
84 performances of nine different algorithms, including RF, ANN, GLM and GBT, to model the current
85 distribution of *Abies alba* and *Fagus sylvatica* for Italy at 1 km resolution using National Forest Inventory
86 (NFI) data; they also evaluated the potential change in the distribution of these species in 2050 according
87 to two of the four Representative Concentration Pathway (RCP) scenarios adopted by IPCC. Prasad et al.
88 (2006) used NFI data as well and were among the first to include topographical and land cover variables to
89 predict the current distribution and future climate scenarios of *Acer saccharum*, *Fagus grandifolia*, *Pinus*
90 *taeda* and *Quercus alba* in the eastern United States using RF and CART at 1 km resolution. Following
91 this trend, Cord et al. (2009) added another layer of complexity: to model the distribution of tree species
92 of the genera *Pinus spp.* and *Quercus spp.* in Mexico, they included multi-temporal remotely sensed data
93 at medium resolution (250 m) as predictors and compared the influence of Earth Observation (EO) data
94 on the predictive performance of the model.

95 Meier et al. (2010) measured how strongly the influence of biotic interactions affects the prediction
96 of species distribution in Swiss broadleaf forests. Biotic interactions are usually included in the form
97 of distribution maps due to data availability: in this study they were included as predictors in the form

98 of relative abundance compared to the target species, relative abundance of large individuals and total
99 shade by large individuals. Results showed how biotic variables considerably improved the predictive
100 performances; their effect were independent from abiotic variables and helped in identifying aspects not
101 easily captured by abiotic predictors only.

102 **Brus et al. (2012)** mapped 23 species groups over Europe at 1 km resolution using NFI data, bioclimatic
103 and topographical variables. Contrary to most SDM studies, where each species is mapped individually
104 and one map per species is provided, **Brus et al. (2012)** use multinomial regression to assign a probability
105 value in the interval [0–1] to each species, with probabilities adding up to 1. This means that each pixel
106 shows the prevalent species, but not if multiple species may coexist in the same area. Another state-of-the-
107 art data source on forest species distribution in Europe is the European Atlas of Forest Tree species. This
108 is among the largest data sources with information on forest tree species for Europe (**San-Miguel-Ayanz,**
109 **J., de Rigo, D., Caudullo, G., Houston Durrant, T., Mauri, 2016**). It describes in detail the auto-ecology
110 of 76 different forest tree species and provides geographical information on each species in the form of
111 chorological maps, probability of presence maps and maximum habitat suitability maps. While these
112 predictions are useful to ascertain the species composition of European forests, there is now a need for
113 spatio-temporal predictions of potential and realized distribution of forest tree species building on new
114 methods that potentially better suited to deal with the increasing resolution and availability of data for
115 both training and prediction.

116 Spatio-temporal modeling, Earth Observation data and specifically the use of high spatial resolution
117 data have only recently started to be explored for SDM (**Gelfand and Shirota, 2021; Pérez Chaves et al.,**
118 **2018; Hefley and Hooten, 2016**). Commonly, SDM still relies on climatic or bioclimatic factors at
119 coarse spatial resolution (≥ 1 km) while in the temporal dimension long time averages (30–50 years)
120 are commonly used (**Iturbide et al., 2018**). The potential impact of differences in resolution of the input
121 variables on the results is often ignored (**Porfirio et al., 2014**) despite the fact the forest spatial patterns in
122 Europe are often linked to management decisions happening at local scales. Furthermore, previous studies
123 comparing predictive performances at different spatial resolutions mention distribution maps with high
124 spatial resolution (< 100 m) and slightly lower prediction accuracy can still be more useful for invasive
125 species management than coarser (> 250 m) but more accurate maps (**Manzoor et al., 2018; Guisan et al.,**
126 **2013; Gottschalk et al., 2011; Prates-Clark et al., 2008**).

127 In the last decade, ecologists have conducted hundreds of studies to determine which methods best
128 suit the needs of SDM. Some general findings are that presence-absence approaches are usually better
129 than presence-only and better predictive power comes at the cost of transferability of the model (**Valavi**
130 **et al., 2021; Pecchi et al., 2019; Guisan et al., 2017**) but no consensus has been reached yet. For example,
131 previous studies have shown that distribution maps derived from the same dataset but using different
132 models can lead to quite opposite conclusions (**Araújo and New, 2007; Pearson et al., 2006**). Inter-model
133 variability in projections has been tackled using ensemble modeling, where numerous independent models
134 are fit using a range of methods applied to the same input data while the outputs of the individual models
135 are aggregated into the final prediction. **Hao et al. (2019)** provide a good general review in this sense,
136 focusing only on SDM studies that use ensemble methods: they reviewed a total of 224 papers, limited,
137 however, to ensemble methods implemented in the BIOMOD software. Strong conclusions could not

138 be drawn due to lack of information on performances of ensemble models versus individual models. On
139 top of that, only few ensemble strategies have been investigated: mean, median and weighted average.
140 These strategies are intuitively simple and involve, in the first two cases, just taking the mean or median
141 of the predictions of the individual models as a final prediction. The weighted average does the same
142 but scales the predictions by weights assigned based on cross validation predictive performances of the
143 models. A robust ensemble technique that, to our knowledge, has not been tested yet for SDM is stacking
144 or stacked generalization. Outputs made by the individual models are the input of a meta-learner which
145 then produces the final prediction (Wolpert, 1992). This approach will be used in this study.

146 To enable more insight into tree species dynamics over EU we have set the following objectives:

- 147 1. To develop a framework for modeling species distribution in space-time with state-of-the-art
148 Machine Learning.
- 149 2. To assess the importance of various sources of Earth Observation data for mapping tree species
150 distributions.
- 151 3. To explore and quantify the importance of high resolution data on model predictive performances.

152 We first review the theory for space-time modeling of forest tree species using presence and absence
153 data. We then show results of our experimental design to build a spatio-temporal ensemble model for
154 SDM based on ML algorithms, variable importance and accuracy assessment for selection of tree species
155 (16) based on 5-fold spatial cross-validation. Finally, we discuss variable importance of various EO data
156 sources and visualize final predictions for sample areas and whole of EU.

157 2 MATERIALS AND METHODS

158 2.1 Species distribution model theory

159 SDM, as any modeling tool, tries to provide a simplified and understandable conceptual representation
160 of a complex phenomenon in the real world, in this case the geographical distribution of a species; it is
161 therefore important to clarify, before any modeling attempt, what kind of distribution a particular study
162 aims to quantify, i.e. what is the ecological meaning of the response variable predicted. Conceptually
163 speaking, the applied procedure is the same: based on the similarity between the values of the predictors
164 associated with the occurrence data, build a function that can be used to assign probability of presence of
165 the target species to areas with unavailable field data. However, based on what is actually modelled, the
166 choice of predictors, training data (presence-only or presence-absence) and modeling techniques may
167 vary (Soberón, 2010).

168 Hutchinson (1957) defined the distribution of a species as the n -dimensional hypervolume in a vectorial
169 space where each component is a variable that positively influences the growth of the species and on
170 which the species has no influence. These variables are usually abiotic/environmental variables, such as
171 temperature, precipitation or elevation (Soberón, 2010; Soberón and Peterson, 2005). By considering the
172 species physiological response to the combined effect of all these variables in the vectorial space, it is
173 possible to capture the fundamental niche of a species, i.e. all the possible conditions that a species can

174 occupy in the environment. The realized niche is instead defined as a narrower region of this vectorial
175 space constrained by biotic interactions and dispersal. The relationship between abiotic factors, biotic
176 interactions and the regions accessible by the species through dispersion was formalized by **Soberón and**
177 **Peterson (2005)** in the BAM diagram.

178 By definition, it is not possible for SDMs to capture the entire fundamental niche because occurrence
179 data gathered through field campaigns is already constrained by biotic interactions. While it is possible to
180 analyze the fundamental niche of a species using mechanistic approaches, SDMs can only investigate the
181 realm of the realized niche (**Guisan and Thuiller, 2005**). On top of the realized niche of a species, in this
182 study we also model the potential niche. The potential niche is defined as that portion of the fundamental
183 niche that exists in the study area at the time of the study (**Peterson et al., 2011**). The concept is relevant
184 as it assesses whether the conditions as defined in the fundamental niche actually exist (**Guisan et al.,**
185 **2017**). **Guisan et al. (2014)** suggested that the traditional approach of SDM to capture just a snapshot
186 of the realized niche may be too restrictive. In the context of climate change, assessing the impact of
187 changing conditions on the distribution of a species requires making predictions in space and time. For
188 space-time modeling, the realized niche has to be built simultaneously considering: (a) different time
189 periods, and (b) different geographical areas (e.g. native area or exotic areas). (**Martinez-Minaya et al.,**
190 **2018; Miller et al., 2007**). The differences in definitions here discussed are then reflected in the way the
191 two niches are modelled:

- 192 1. Potential distribution is modelled using variables that capture the ecological gradients of the study
193 area and that are in no way influenced by other species (i.e. no biotic interactions),
- 194 2. Remotely sensed, reflectance-based data (i.e. land cover, spectral indices etc.) are not used to
195 model potential distribution as vegetation species distribution may be influenced owing to human
196 alteration of the landscape, underestimating the potential distribution, or capturing disturbed areas
197 that are just temporary not suitable for the species (**Bradley et al., 2012**). The same is true for other
198 human altered of environmental factors such as the Urban Heat Island (UHI) effect (**Bechtel and**
199 **Schmidt, 2011; Pigeon et al., 2007**),
- 200 3. The realized distribution is modelled using a combination of environmental variables (temperature,
201 precipitation, topography), biotic variables and reflectance-based variables.

202 Occurrence data used in SDM is either presence-only (i.e. records of locations where the species was
203 observed) or presence-absence (i.e. locations of both presence and absence of a species have been recorded
204 in the sampling design): ML algorithms for SDM treat the suitability task as a binary classification
205 problem, where the response variable can assume a discrete value of 1 or 0 (i.e. binary response for
206 respectively, presence or absence) or a continuous value from 0 to 100 (i.e. probability of presence), with
207 the latter being the most used (**Guisan et al., 2017**). Species occurrence datasets are hard to obtain due
208 to the high costs of the field surveys. Presence and true absence data usually suffer from preferential
209 sampling, meaning that a) selection of the sampling locations is not independent from the response
210 variable (i.e. species occurrence) and b) the locations in the study area are not equally likely to be sampled
211 (**Martinez-Minaya et al., 2018**). True absence data is usually not available but replaced with background

212 data (Phillips et al., 2009), which does not try to infer where the species is absent but rather to characterize
213 the environmental background of the study area. Alternatively, absence data is artificially generated and
214 referred to as pseudo-absence. Generally, pseudo-absence indicates locations that researchers assume are
215 potentially unsuitable for the target species, and points can be generated in several ways (Iturbide et al.,
216 2015). Different kinds of absence data have to be used depending on the distribution modelled: potential
217 distribution cannot be modelled with absence data coming from non environmental factors (i.e. a species
218 being absent due to competition/biotic interactions) (Jiménez-Valverde et al., 2008).

219 As a source of homogeneously distributed true absence data for both potential and realized distribution,
220 in this study we will use the Land Use/Cover Area Survey (LUCAS) (EUROSTAT, 2017) dataset: in-situ
221 observations of land use and land cover distributed on a 2 x 2 km grid covering the whole European
222 Union (see d’Andrimont et al. (2021) for more information and [https://ec.europa.eu/eurostat/
223 web/lucas/data/lucas-grid](https://ec.europa.eu/eurostat/web/lucas/data/lucas-grid) for the official grid).

224 2.2 General workflow

225 Potential and realized distribution maps were produced for continental Europe for the time period January
226 2000 – December 2020 using a spatio-temporal modeling approach. Final prediction maps show the
227 probability of presence (0–100%) of at least one individual of the target species in the area covered by a
228 30 m pixel. Probability of presence is relative to the mapped target species, irrespective of the potential
229 co-occurrence of other species in the same 30 m pixel and should not be confused with the absolute
230 abundance or proportion of each species in the pixel area. The sum of the presence probabilities of
231 different species in the same pixel can thus exceed 100 %. With each probability map comes a model
232 uncertainty map, where the user can check the reliability of probability map values per pixel. We produced
233 one potential distribution map and six realized distribution maps for each species: the assumption is that
234 the conditions in the study area that determine the potential distribution of the species did not change
235 over the time period analyzed; this does not hold for the realized distribution. We split the time period
236 analyzed in six time windows according to the following scheme:

- 237 • 2000–2002,
- 238 • 2002–2006,
- 239 • 2006–2010,
- 240 • 2010–2014,
- 241 • 2014–2018,
- 242 • 2018–2020.

243 One realized distribution map was produced for each time period. The general workflow used to derive the
244 distribution maps is shown in Fig. 1: two datasets for each species (potential and realized dataset) were
245 produced and overlaid with a set of both static and dynamic predictor variables, matching both the point
246 location in space (x, y) and time (t). Predictor variables used for potential and realized distribution follow
247 the rules listed in the previous section: reflectance-based predictor variables and species distribution maps

248 were used to model only the realized distribution. The classification matrix produced was used to fit a
 249 spatio-temporal ensemble model based on stacking. The workflow was implemented in the Python and R
 250 programming languages.

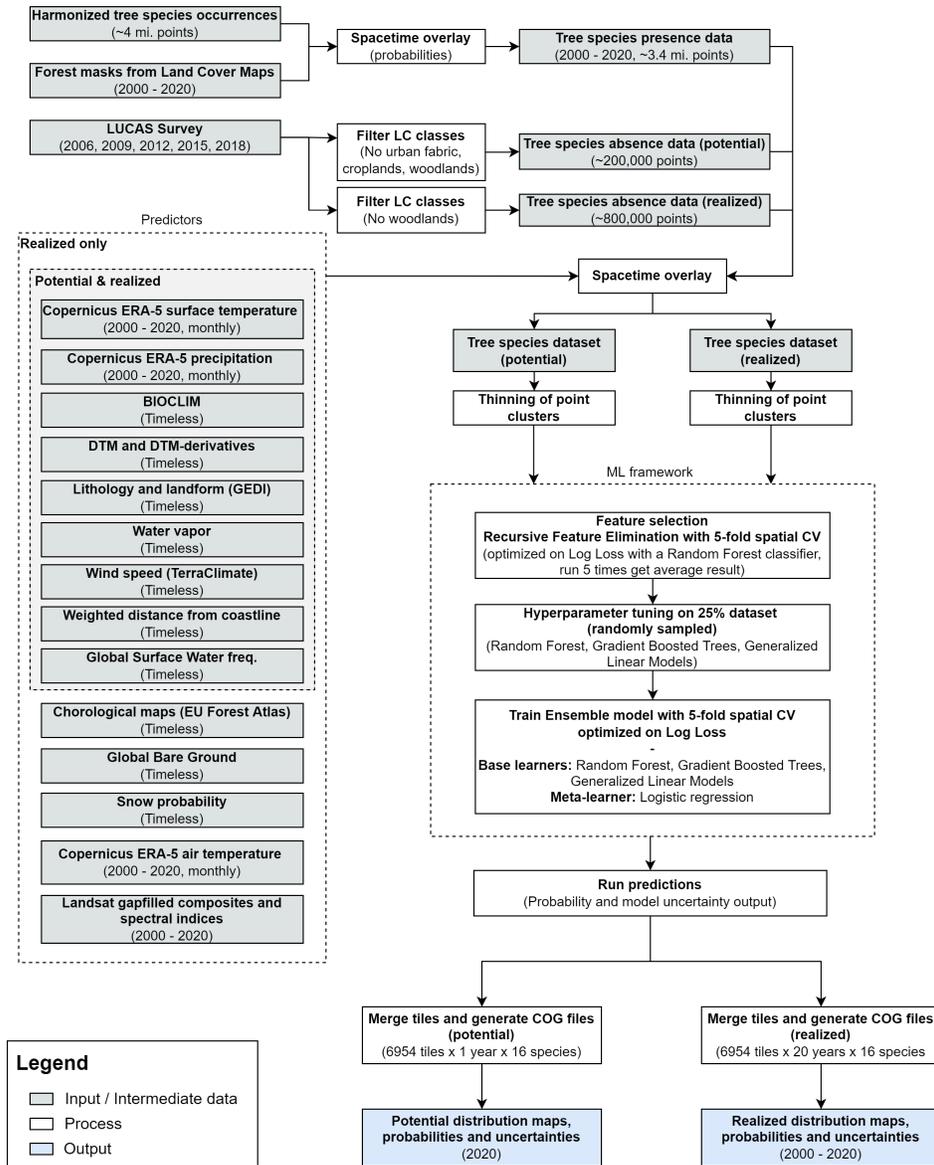


Figure 1. General workflow illustrating the preparation of the point data, the predictor variables used, model building (feature selection — hyperparameter optimization — training) and preparation of distribution maps.

251 2.3 Study area

252 The study area covers the European continent, that is all countries included in the Corine Land Cover
 253 (CLC) database (Büttner et al., 1998) except Turkey. European forests cover 33% of the continent's land
 254 area. Owing to the variety of climatic conditions across both latitudinal and longitudinal gradients, twelve
 255 out of the 20 FAO Forest Ecological Zones are represented in European forests (de Rigo et al., 2016b).

256 The European Atlas of Forest Tree Species (San-Miguel-Ayanz, J., de Rigo, D., Caudullo, G., Houston
257 Durrant, T., Mauri, 2016) reports detailed information for a total of 76 forest tree species. From those, the
258 following 16 were chosen and modelled in this study:

- 259 1. Silver fir (*Abies alba* Mill.),
- 260 2. Sweets chestnut (*Castanea sativa* Mill.),
- 261 3. Common hazel (*Corylus avellana* L.),
- 262 4. European beech (*Fagus sylvatica* L.),
- 263 5. Olive tree (*Olea europaea* L.),
- 264 6. Norway spruce (*Picea abies* L. H. Karst.),
- 265 7. Aleppo pine (*Pinus halepensis* Mill.),
- 266 8. Austrian pine (*Pinus nigra* J. F. Arnold),
- 267 9. Stone pine (*Pinus pinea* L.),
- 268 10. Scots pine (*Pinus sylvestris* L.),
- 269 11. Sweet cherry (*Prunus avium* L.),
- 270 12. Turkey oak (*Quercus cerris* L.),
- 271 13. Holm oak (*Quercus ilex* L.),
- 272 14. Common oak (*Quercus robur* L.),
- 273 15. Cork oak (*Quercus suber* L.),
- 274 16. Goat willow (*Salix caprea* L.).

275 **2.4 Training points**

276 **2.4.1 Preparing and combining legacy occurrence points**

277 A total of 2,454,997 tree species occurrence locations from three different sources was gathered. The
278 majority of points (71%) comes from the Global Biodiversity Information Facility (GBIF) which is an
279 open database fed by field observations from individual researchers. National forest inventory data from
280 multiple EU member states published by Mauri et al. (2017) forms another 23% of the dataset. The
281 remaining 6% comes from the LUCAS dataset.

282 All sources were filtered, harmonized and supported by additional information. Entries were filtered
283 for species included in the European Atlas of Forest Tree Species (San-Miguel-Ayanz, J., de Rigo, D.,
284 Caudullo, G., Houston Durrant, T., Mauri, 2016). Occurrences with a taxonomy rank other than species
285 or genus were omitted. Same applies to points which had flags indicating serious location issues (i.e.
286 missing coordinates). Geometries were re-projected to coordinate reference system ETRS89 / LAEA
287 Europe (EPSG: 3035). A high resolution land mask for Europe (Hengl et al., 2020) was applied to further

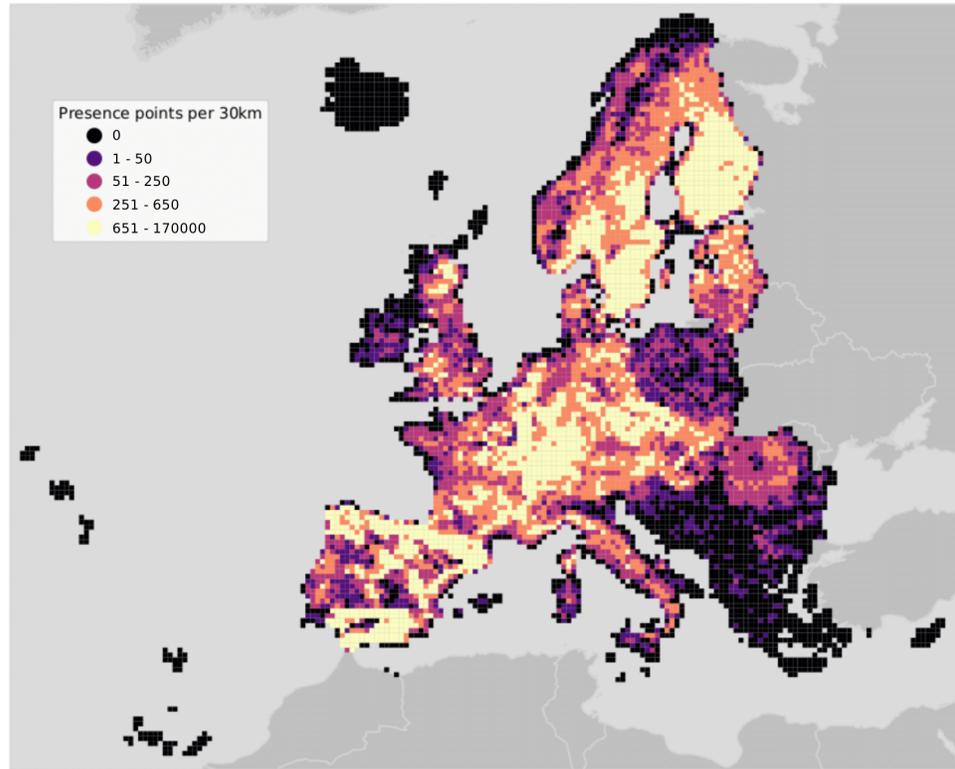


Figure 2. Map of the study area overlaid with a grid of 30 km tiles used for spatial 5-fold cross validation. Pixel value shows the number of presence points per tile.

288 exclude misplaced occurrence points. GBIF taxon and genus keys were derived for the other two data
 289 sources. Quality flag variables for location accuracy and date were established from existing metadata
 290 to indicate potentially problematic entries. The harmonized point dataset has information on species
 291 and genus (including respective GBIF keys), year of observation, country, original data source, citation,
 292 and license among other auxiliary variables. Its generation is reproducible. The dataset was published
 293 separately and is openly available along with the code and individual GBIF dataset citations (Heisig and
 294 Hengl, 2020).

295 We used yearly forest masks derived from Witjes et al. (2021) to decide upon including point data
 296 lacking the year of observation. Witjes et al. (2021) provides yearly probability maps at 30 m for the
 297 period 2000–2020 for 43 land cover classes according to the CLC level 3 legend. We overlaid the points
 298 with the probability maps for the classes:

- 299 • 311: Broad-leaved forest,
- 300 • 312: Coniferous forest,
- 301 • 313: Mixed forest,
- 302 • 323: Sclerophyllous forest,

- 303 • 324: Transitional woodland-shrub,
- 304 • 333: Sparsely vegetated area.

305 Points were used only if the probability value extracted for at least one of the classes was $\geq 50\%$ for
306 all the years considered. Each unique combination of longitude, latitude and year was then considered
307 as an independent sample. An additional quality flag was added to distinguish points coming from this
308 operation and the points with original year of observation coming from source datasets.

309 **2.4.2 Preparing non-occurrence points**

310 A total of 883,630 land cover points was gathered from the LUCAS database as provided by Eurostat
311 and used as absence data. All LUCAS survey data (2006, 2009, 2012, 2015 and 2019) was used: each
312 survey was first downloaded individually and then aggregated. As for the occurrence points, each unique
313 combination of longitude, latitude and year was considered as an independent sample. The survey assigns
314 each location as belonging to one of the following 8 main categories:

- 315 • A: Artificial land,
- 316 • B: Cropland,
- 317 • C: Woodland,
- 318 • D: Shrubland,
- 319 • E: Grassland,
- 320 • F: Bareland,
- 321 • G: Water,
- 322 • H: Wetlands.

323 The above classes were used for selecting observations for the absence dataset. Points in class C were
324 excluded as absence data as that class already served the selection of presence data. For modelling
325 the actual distribution, all remaining points were included in the first selection, while for the potential
326 distribution points coming from classes A and B were excluded. For modelling the potential distribution
327 the selection of absence points was more restrictive. Points were overlaid with a rasterized chorological
328 map for each of the target species downloaded from the European Atlas of Forest Tree Species portal
329 (see subsection: 2.1). Only points falling outside the area indicated by the chorological map were used as
330 absence for the potential distribution.

331 **2.5 Predictor variables**

332 A total of 585 harmonized variables covering continental Europe (Hengl et al., 2020) at different spatial
333 resolution were used as predictors to model the realized distribution of the species. In this study we
334 included both dynamic (i.e. time-series of data of different temporal resolution) variables covering the
335 time period January 2000 – December 2020 and static (i.e. variables not expected to change during the

336 modelled time period) variables. A subset of only 318 variables were used instead to model the potential
 337 distribution (see Fig. 1). All data was reprojected in the Projected coordinate system for Europe (EPSG
 338 code: 3035) before the analysis.

339 2.5.1 Dynamic data

340 We used a reprocessed version of Landsat ARD data provided by GLAD (Potapov et al., 2020): time
 341 series used in this study covers the period 1999–2020. Cloud and cloud shadow pixels were removed
 342 from the images, maintaining only the quality assessment-QA values labeled as clear-sky. Afterwards,
 343 individual images were averaged by season according to three different quantiles (25th, 50th and 75th)
 344 and the following calendar dates for all period:

- 345 • Winter: December 2 of previous year until March 20 of current year,
- 346 • Spring: March 21 until June 24 of current year,
- 347 • Summer: June 25 until September 12 of current year,
- 348 • Fall: September 13 until December 1 of current year.

349 84 images (3 quantiles \times 4 seasons \times 7 Landsat bands) were produced for each year. Missing values were
 350 imputed using the *Temporal Moving Window Median* algorithm. For more details on the preprocessing of
 351 Landsat data for this study see Witjes et al. (2021). 7 different spectral indices were computed for each
 352 year and season using the 50th quantile only, for a total of $7 \times 4 = 28$ spectral indices variables per year.

Table 1. Table with Landsat-derived spectral indices used in this study.

Spectral Index	Abbreviation	Formula	Reference
Enhanced Vegetation Index	EVI	$2.5 \times \frac{NIR - RED}{NIR + 6 \times RED - 7.5 \times BLUE + 1}$	(Huete et al., 2002)
Enhanced Vegetation Index 2	EVI2	$2.5 \times \frac{NIR - RED}{NIR + 2.4 \times RED + 1}$	(Jiang et al., 2008)
Modified Soil Adjusted Vegetation Index	MSAVI	$\frac{(2 \times NIR + 1) - \sqrt{(2 \times NIR + 1)^2 - 8 \times (NIR - RED)}}{2}$	(Qi et al., 1994)
Normalized Burned Ratio	NBR	$\frac{NIR - SWIR2}{NIR + SWIR2}$	(Key and Benson, 1999)
Normalized Difference Vegetation Index	NDVI	$\frac{NIR - RED}{NIR + RED}$	(Tucker, 1979)
Normalized Difference Wetness Index	NDWI	$\frac{NIR - SWIR1}{NIR + SWIR1}$	(Gao, 1996)
Soil Adjusted Vegetation Index	SAVI	$(1 + 0.5) \times \frac{NIR - RED}{(NIR + RED + 0.5)}$	(Huete, 1988)

353 A reprocessing of the ERA5 Land hourly dataset has been used to have monthly aggregates of air
 354 temperature (2 meters above ground), surface temperature and precipitation. Original ERA5 data was
 355 aggregated to daily data, and subsequently to monthly data, with increased resolution (1 km) using
 356 CHLSA data (Karger et al., 2020): in this way the general spatial and temporal pattern of ERA5 Land
 357 dataset was kept while using the fine spatial detail coming from the CHLSA dataset. The following
 358 steps were used for temperature data:

- 359 1. aggregate CHLSA to ERA5 spatial resolution,

- 360 2. calculate difference between ERA5 Land and aggregated CHELSA,
- 361 3. interpolate differences with a Gaussian filter to 30 arc seconds,
- 362 4. add the interpolated differences to CHELSA.

363 A different approach was used for precipitation, with proportions instead of differences: using proportions
364 ensures that areas without recorded precipitation remain areas without precipitation; only in the case
365 of actual precipitation in a given area, precipitation was redistributed according to the spatial detail of
366 CHELSA:

- 367 1. aggregate CHELSA to ERA5 spatial resolution,
- 368 2. calculate proportion between ERA5 Land and aggregated CHELSA,
- 369 3. interpolate proportion with a Gaussian filter to 30 arc seconds,
- 370 4. multiply the interpolated proportion with CHELSA.

371 For air and surface temperature we obtained the monthly minimum, mean and maximum, while for
372 precipitation the monthly sum. We then computed the standard deviation of each of these variables
373 and added these as additional predictor variables for modeling. To simulate the cumulative effect of
374 temperature and precipitation regimes on a short term scale (5 years), we also calculated the monthly
375 average of the previous 5 years for each variable, for a total of 336 climatic time series layers.

376 2.5.2 Static covariate datasets

377 As additional static covariates, we used the following datasets:

- 378 • 19 bioclimatic variables (Hijmans et al., 2005) for the period 1979 - 2013 to provide a baseline
379 of the actual state of the climate; we however decided use the bioclimatic variables from the
380 CHELSA dataset (1 km spatial resolution) since CHELSA has been claimed to better match data
381 from meteorological stations than WorldClim (Karger et al., 2017). Layers were freely downloaded
382 from <https://chelsa-climate.org/bioclim/>,
- 383 • 50 different chorological maps downloaded from the Europa Atlas of Forest Tree Species web
384 portal. Chorological maps provide a qualitative overview of the spatial distribution of a species
385 over an area, differentiating between native and introduced. We considered both the native and
386 introduced areas as the potential distribution of a species for the time period covered by the study.
387 The maps are harmonized products derived from different historical bibliographic sources or from
388 the EUFORGEN website (<http://www.euforgen.org>); they are provided as vectorial layers
389 and the minimum mapping unit may vary across the species. To include these layers in the models,
390 we first rasterized the maps on a 10 km grid, assigning a value of 100 or 0, respectively, where
391 the species could potentially be present or not. We then resampled the maps at a spatial resolution
392 of 30 m using the cubic spline method to create a smoothed, continuous area around the reported
393 range,

- 394 • Global bare ground cover from Hansen et al. (2013). The layer is derived from annual com-
395 posites from 2000 to 2012 and provides information on bare ground cover on a percent (1–
396 100) scale at 30 m resolution; it was downloaded from [https://glad.umd.edu/dataset/
397 global-2010-bare-ground-30-m](https://glad.umd.edu/dataset/global-2010-bare-ground-30-m),
- 398 • Solar direct and diffuse irradiation images at 1 km resolution were downloaded from [https://
399 //globalsolaratlas.info/download](https://globalsolaratlas.info/download),
- 400 • 13 cloud fraction layers (monthly averages and annual average) derived from MODIS (Wilson and
401 Jetz, 2016) at 1 km resolution obtained from <https://www.earthenv.org/cloud>,
- 402 • Digital terrain model (DTM) for Europe (Hengl et al., 2020) and DTM-derived (slope, hillshade)
403 variables at 30 m resolution, available along with its reproducible code at [https://doi.org/10.
404 5281/zenodo.4724549](https://doi.org/10.5281/zenodo.4724549),
- 405 • Easternness, northness derived in GRASS GIS (Olaya, 2009), and positive and negative openness
406 derived using SAGA GIS and available via <https://doi.org/10.5281/zenodo.4486135>,
- 407 • Probability of surface water occurrence at 30 m resolution derived from Landsat time series (Pekel
408 et al., 2016), obtained from and freely available at [https://global-surface-water.appspot.
409 com/](https://global-surface-water.appspot.com/),
- 410 • Height above nearest drainage (HAND) and flow accumulation area at 90 m resolution from the
411 MERIT Hydro global hydrography datasets, freely available at [http://hydro.iis.u-tokyo.
412 ac.jp/~yamada/MERIT_Hydro/](http://hydro.iis.u-tokyo.ac.jp/~yamada/MERIT_Hydro/),
- 413 • Long-term flood hazard map calculated on a 500 years time period at 1 km resolution (Dottori et al.,
414 2016),
- 415 • Continental Europe surface lithology based on the European Geological Data Infrastructure (EGDI)
416 at 1:1M scale produced by GEOZS, Slovenia, and available at [https://doi.org/10.5281/
417 zenodo.4787631](https://doi.org/10.5281/zenodo.4787631),
- 418 • Water vapor pressure (kPa) based on the WorldClim2.1 dataset (Fick and Hijmans, 2017); obtained
419 from <http://www.worldclim.com/version2>,
- 420 • Long-term snow probability (2000–2012) at 500 m resolution based on the MODIS (MOD10A2)
421 and available at <https://doi.org/10.5281/zenodo.5774953>,
- 422 • Monthly wind speed (1998–2018) at 5 km resolution from TerraClimate available at [https://
423 //www.climatologylab.org/terraclimate.html](https://www.climatologylab.org/terraclimate.html).

424 2.6 Species datasets

425 Points were overlaid with the predictor variables and two presence-absence datasets were produced for
426 each species, one to be used for potential distribution and one for realized distribution. The final dataset
427 contains predictor variables plus 8 columns with metadata used to uniquely identify the points:

- 428 • id: unique point identifier,
- 429 • year: year of observation,
- 430 • postprocess: quality flag to identify if the temporal reference of an observation comes from the
431 original dataset or is the result of spatiotemporal overlay with forest masks,
- 432 • tile id: contains the tile id from the 30 km grid,
- 433 • easting: longitude coordinates in Coordinate Reference System ETRS89 / LAEA Europe (= EPSG
434 code 3035),
- 435 • northing: latitude coordinates in Coordinate Reference System ETRS89 / LAEA Europe (= EPSG
436 code 3035),
- 437 • atlas class: name of the tree species according to the European Atlas of Forest Tree Species or
438 NULL in case of absence point,
- 439 • lc1: contains original LUCAS land cover class or NULL if it's a presence point.

440 While absence points are homogeneously distributed in the study area, the same is not true for presence
441 points. Fig. 2 shows the number of presence points for all the species selected in this study aggregated
442 using a 30 km grid: clusters of points can be observed in Scandinavia, France and Spain, while some
443 countries (i.e. Cyprus, Iceland) have no points at all.

444 To obtain a spatially balanced presence-absence dataset for each species, we empirically defined a
445 fixed amount of points to be selected from each tile: we first overlaid the points with a 30 km grid and
446 counted the amount of points per tile. We selected the highest number scored by each species and then
447 calculated the median of the distribution made by all the highest scores. Results of this thinning operation
448 can be seen in Table 2 and Fig. 3.

449 2.7 Feature selection

450 Features for potential and realized distribution for each species were selected using the Recursive Feature
451 Elimination (RFE) strategy, implemented in the Python programming language and publicly available in
452 the scikit-learn library. For each combination of species and modelled distribution we trained a random
453 forest classifier (num.trees = 50, default values were used for the other parameters): RFE fits the model
454 and removes the weakest feature (or features) until a specified number of features is reached, then ranks
455 the importance of the features based on the model's coefficients (for regression-based models) or feature
456 importance (for random forest).

457 The minimum number of features was not known before hand: to select this number, we ran the
458 Recursive Feature Elimination with a spatial 5-fold Cross Validation (RFECV), using the logarithmic loss,
459 or logloss, as a scoring estimator. Logloss is one of the most robust performance metric when it comes to
460 imbalanced datasets (Ferri et al., 2009). Logloss is indicative of how close the predicted probability for

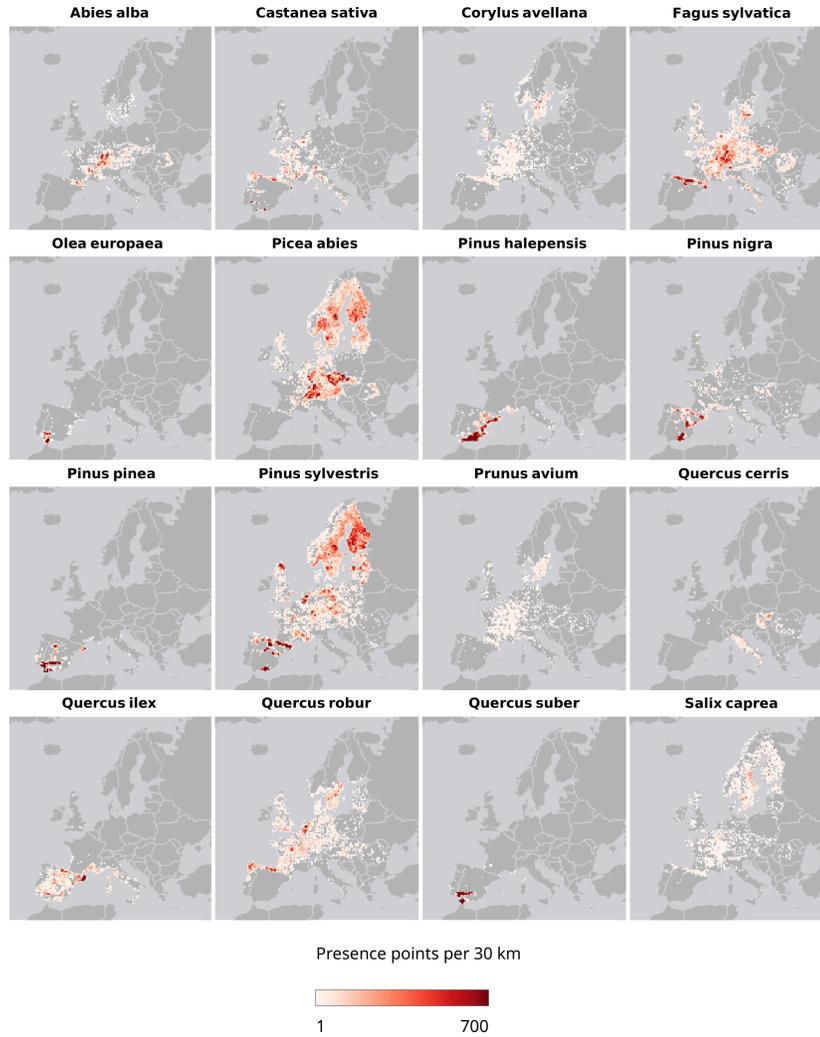


Figure 3. Distribution of presence points per species after thinning. Absence points are omitted for visualization purposes.

461 an observation i is to the corresponding label y . For binary classification with label $y \in \{0, 1\}$ the overall
 462 logloss was calculated as:

$$f_{\text{Logistic}}^* = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \ln [p(y_i)] + (1 - y_i) \cdot \ln [1 - p(y_i)] \quad (1)$$

463 where N is the total number of observations and $p(y_i)$ is the predicted probability for an observation with
 464 $y_i = 1$. It follows that values close to 0 indicate high prediction performances, with logloss = 0 being a
 465 perfect match, and values that are positive to infinite are progressively worse scores. For comparison, the
 466 value of logloss for random assignment depends on the number of classes (a) and the prevalence of the
 467 classes (b): for binary classification and a balanced (50:50) dataset with $N = 10$ observations, the equation
 468 (1) gives a value of 0.69.

469 We ran the RFECV on a 25% random subsample for each species and modelled distribution; this

Table 2. Number of presence and absence points used to model potential and realized distribution for each species before and after the thinning operation.

Distribution	Species	Presence	Absence	Presence thinned	Absence thinned	Prevalence	Prevalence thinned
Potential	<i>Abies alba</i>	45,495	108,603	45,410	108,445	0.42	0.42
Potential	<i>Castanea sativa</i>	77,382	107,667	49,269	107,363	0.72	0.46
Potential	<i>Corylus avellana</i>	32,321	64,135	32,141	63,955	0.50	0.51
Potential	<i>Fagus sylvatica</i>	197,240	55,595	180,817	55,266	3.50	3.30
Potential	<i>Olea europaea</i>	50,656	191,830	12,761	191,769	0.26	0.07
Potential	<i>Picea abies</i>	360,271	172,879	351,063	172,769	2.10	2.00
Potential	<i>Pinus halepensis</i>	233,964	180,795	71,258	179,931	1.30	0.40
Potential	<i>Pinus nigra</i>	139,516	225,382	54,739	224,752	0.62	0.24
Potential	<i>Pinus pinea</i>	239,254	221,263	41,188	221,050	1.10	0.19
Potential	<i>Pinus sylvestris</i>	507,681	44,652	415,859	44,077	11.00	9.40
Potential	<i>Prunus avium</i>	22,978	88,517	22,862	88,401	0.26	0.26
Potential	<i>Quercus cerris</i>	13,856	104,782	13,774	104,700	0.13	0.13
Potential	<i>Quercus ilex</i>	57,690	203,273	52,213	203,207	0.28	0.26
Potential	<i>Quercus robur</i>	113,044	52,493	111,399	52,190	2.20	2.10
Potential	<i>Quercus suber</i>	419,975	211,403	26,361	211,253	2.00	0.12
Potential	<i>Salix caprea</i>	45,967	76,263	45,789	76,082	0.60	0.60
Realized	<i>Abies alba</i>	45,495	558,564	44,882	558,493	0.08	0.08
Realized	<i>Castanea sativa</i>	77,382	558,564	48,969	554,503	0.14	0.09
Realized	<i>Corylus avellana</i>	32,321	558,564	32,141	554,781	0.06	0.06
Realized	<i>Fagus sylvatica</i>	197,240	558,564	179,108	557,490	0.35	0.32
Realized	<i>Olea europaea</i>	50,656	558,564	12,483	542,354	0.09	0.02
Realized	<i>Picea abies</i>	360,271	558,564	348,682	557,571	0.64	0.63
Realized	<i>Pinus halepensis</i>	233,964	558,564	68,400	553,831	0.42	0.12
Realized	<i>Pinus nigra</i>	139,516	558,564	53,717	557,050	0.25	0.10
Realized	<i>Pinus pinea</i>	239,254	558,564	39,766	555,684	0.43	0.07
Realized	<i>Pinus sylvestris</i>	507,681	558,564	412,502	555,730	0.91	0.74
Realized	<i>Prunus avium</i>	22,978	558,564	22,862	558,556	0.04	0.04
Realized	<i>Quercus cerris</i>	13,856	558,564	13,774	558,559	0.03	0.03
Realized	<i>Quercus ilex</i>	57,690	558,564	52,213	558,280	0.10	0.09
Realized	<i>Quercus robur</i>	113,044	558,564	110,748	558,267	0.20	0.20
Realized	<i>Quercus suber</i>	419,975	558,564	25,617	556,598	0.75	0.05
Realized	<i>Salix caprea</i>	45,967	558,564	45,769	558,558	0.08	0.08

470 operation was replicated 5 times. For each iteration we selected the minimum of the function showed in
 471 Fig. 4 and the averaged result was then used as the minimum number of features for the RFE.

472 2.8 Model building and evaluation

473 2.8.1 Modeling methods

474 To build an ensemble model, we decided to compare predictive performances and computing time
 475 (hyperparameter tuning — cross validation — prediction time) of different machine learning algorithms
 476 on a random 25% subset of observations for both potential and realized distribution datasets. A detailed
 477 workflow of this process is shown in Fig. 5. We decided to conduct this test on seven different species:
 478 choice of the species was based on the spatial distribution of the training points and the ratio between
 479 presence and absence points. In this way, algorithms performances can be tested on different ecological
 480 conditions (latitudinal and longitudinal gradient) and imbalance of classes. The species selected were:
 481 *Abies alba*, *Castanea sativa*, *Fagus sylvatica*, *Picea abies*, *Pinus halepensis* and *Pinus sylvestris*.

482 The following algorithms were compared:

- 483 • Random Forest (RF) (Breiman, 2001a),
- 484 • Gradient-boosted trees (GBT) (Friedman, 2002),
- 485 • Classification trees (CART) (Theureau and Atkinson, 2011),

- 486 • Generalized Linear Models (Nelder and Wedderburn, 1972) with Lasso regularization (Tibshirani,
487 1996) (just GLM from now on),
- 488 • C5.0 (Quinlan, 1986),
- 489 • K-nearest neighbor (KNN) (Fix and Hodges, 1989),
- 490 • Artificial Neural Network (ANN) (Ripley and Venables, 2017),

491 Analyses were conducted using the *mlr* (Machine Learning for R) package (Bischi et al., 2016): the
492 package provides a standardized framework that wrap different machine learning algorithms already
493 implemented in R in other independent packages. It also provides hyperparameter tuning with different
494 optimization strategies, resampling functionalities (bootstrapping, cross validation), benchmarking and
495 visualizations through the *ggplot2* package. For each algorithm, an hyperparameter space was defined:
496 combinations of hyperparameters were generated per model based on a grid search of 5 steps per
497 hyperparameter. Due to computational constraints, we set the *num.trees* parameter for Random forest to
498 85. The tested activation functions for the neural network were *sigmoid* and *tanh*, while for the output we
499 tested both *sigmoid* and *softmax*. For GLM we used the automatically generated λ sequence and selected
500 the λ_{\min} . The rest of the hyperparameter ranges are shown in Table 3.

501 2.8.2 Selecting component models

502 We evaluated each combination of hyperparameters by comparing logarithmic loss values during a 5-fold
503 spatial cross validation replicated 5 times: we used spatial cross validation instead of normal cross
504 validation for hyperparameter tuning because it educes overoptimistic performance results in the presence
505 of strong data clustering (Schratz et al., 2019). We used the tile ID produced in the tiling system for

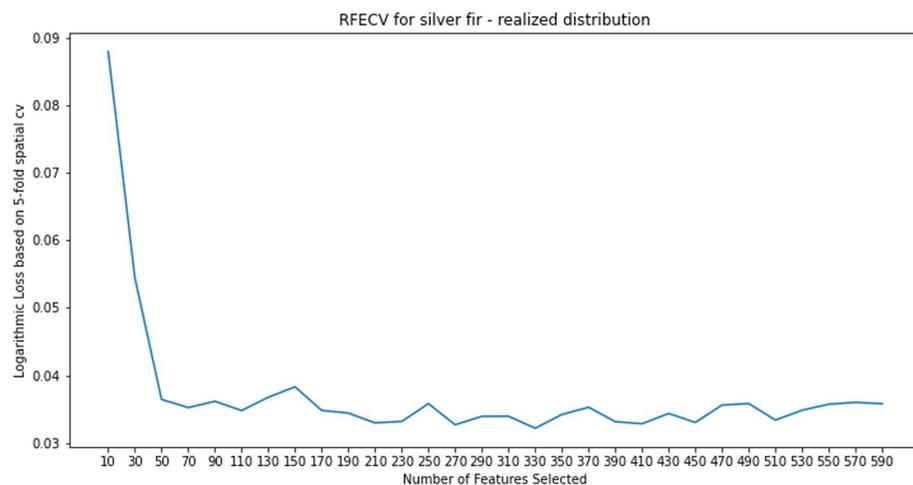


Figure 4. Log Loss performances by number of selected features. In this iteration ran for the realized distribution of the silver fir, 330 is the minimum of the function.

506 Europe (Coordinate Reference System ETRS89 / LAEA Europe, epsg:3035) as the blocking parameter
 507 in the training function in *mlr*. All the compared algorithms were used in "probability" mode, that is,
 508 predicting for each observation in the dataset a probability value for presence (class 1) and absence (class
 509 0). Besides the performance achieved in the logloss metric, computing time for the hyperparameter
 510 tuning, a 5-fold spatial cross validation and prediction time for a 30 km tile were also considered as
 511 additional criteria: we calculated these two metrics only for the species that had the highest computational
 512 costs (*Pinus sylvestris*). This gave us an estimate of how long the process of training each component
 513 model could take during the building of the ensemble model. We used logloss performance as the first
 514 criteria to choose the component models: only in the case of two or more methods performing within one
 515 standard deviation from the average performance, we chose the computationally fastest one across all
 516 areas (hyperparameter tuning — cross validation — prediction) and tasks (potential — realized).

517 2.8.3 Training ensemble model using stacking

518 Stacked generalization involves combining predictions made by level 0 models and using them as training
 519 data for a level 1 model (or meta-model or meta-learner from now on) (Wolpert, 1992). To avoid or limit
 520 overfitting, there are currently two approaches used to prepare the training data: via hold-out method,
 521 which splits the original dataset in training and test, or via k-fold cross validation. We used the latter in
 522 this study, via a 5-fold spatial cross validation: the out-of-fold predictions were used to build a level 1

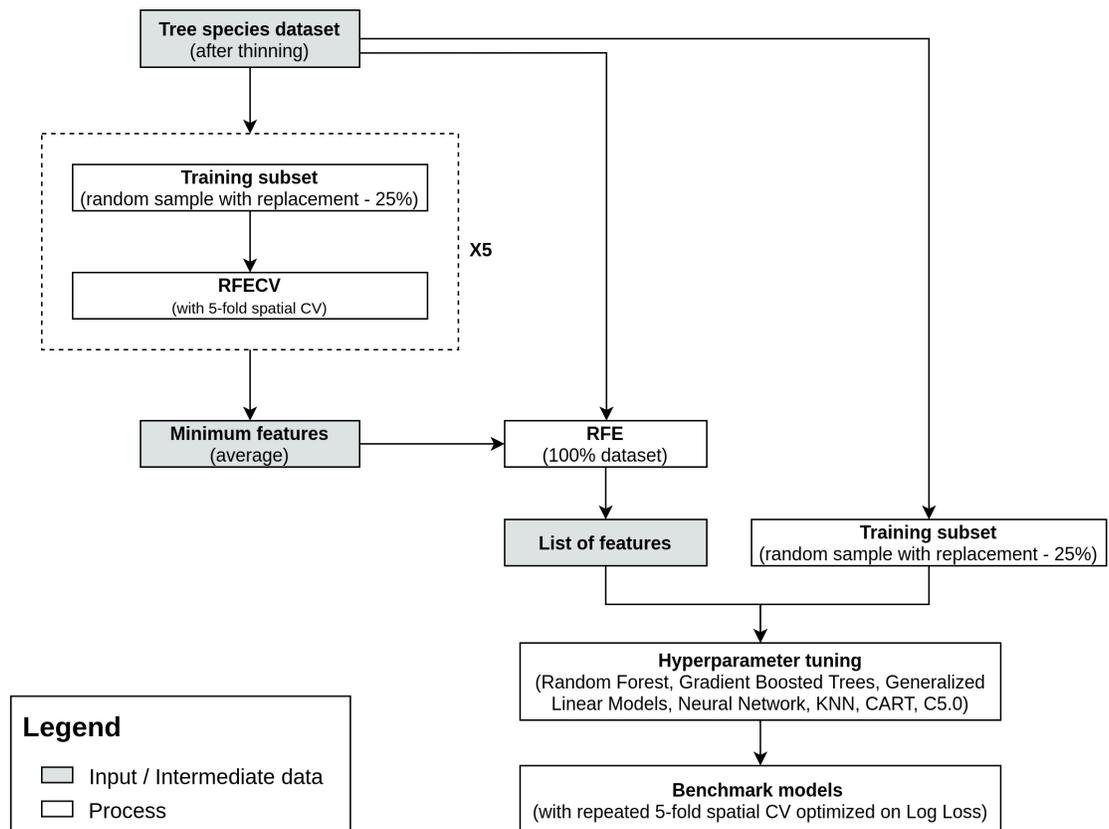


Figure 5. Example workflow illustrating the feature selection and benchmarking process for one species

Table 3. Hyperparameter space for the analyzed algorithms. In light gray the name of the R package used to implement the algorithm is reported, in brackets the name of the algorithm. p refers to the number of predictor variables, while columns *Lower* and *Upper* indicate the bounds of the regions in the hyperparameter space.

Algorithm	Hyperparameter	Type	Lower	Upper
C5.0				
(Classification trees)	minCases	integer	0	10
	CF	numeric	0	0.5
kkn				
(k-nearest neighbor)	k	integer	1	50
deepnet				
(Artificial neural network)	learning rate	numeric	0.0001	0.00001
	numepochs	integer	10	20
	batchsize	integer	50	150
	hidden_dropout	numeric	0.1	0.3
	activationfunction	discrete	-	-
	output	discrete	-	-
	momentum	numeric	0	0.05
	number.of.layers	integer	2	4
	units	integer	32	64
ranger				
(Random forest)	mtry	integer	$\sqrt{p}/3$	p
rpart				
(CART)	minsplit	integer	20	25
	minbucket	integer	5	10
	cp	numeric	0.01	0.1
	maxcompete	integer	3	4
	maxsurrogate	integer	4	5
	usesurrogate	discrete	-	-
	surrogatestyle	discrete	-	-
	maxdepth	integer	5	15
xgboost				
(gradient-boosted trees)	nrounds	integer	10	20
	max_depth	integer	3	5
	eta	numeric	0.01	0.1
	subsample	numeric	0.5	0.9
	min_child_weight	integer	10	20
	colsample_bytree	numeric	0.5	0.9

523 training dataset for the meta-learner. We used logistic regression with Lasso regularization (Tibshirani,
524 1996) as a meta-learner, which is usually the most used model for classification problems (Gomes et al.,
525 2012). Final predictions are delivered as probability maps (0–100%) for presence together with model
526 uncertainty maps: we consider as model uncertainty the standard deviation of the predicted values of
527 the base learners. The principle is that the higher the standard deviation the more uncertain the model is
528 towards the right value to assign to the pixel (Brown et al., 2020).

529 2.8.4 Variable importance assessment

530 To assess to what extent the three level 0 models used different parts of the available feature space and
531 the agreement between these models, we compared the variable importance when possible. For Random
532 forest and CART we used Gini importance, for C5.0 the "percentage of training set samples that fall into

533 *all the terminal nodes after the split*”(Quinlan, 1986), for Gradient boosted trees the gain metric (Shi
 534 et al., 2019) and for Generalized linear models the coefficients for the minimum fitted value of λ (Hastie
 535 et al., 2016).

536 2.8.5 Model evaluation

537 Predictive performances of the ensemble model was assessed through spatial 5-fold cross-validation
 538 repeated 5 times (Roberts et al., 2017) with logloss as performance metric. To investigate if the ensemble
 539 model outperformed the component models, we compared results of the spatial cross validation of the
 540 ensemble with the results of the component models. The area under the ROC curve (AUC) is a commonly
 541 used metric to evaluate SDMs predictive performances due to being threshold independent (Shabani et al.,
 542 2018); however, it is also more sensitive to prevalence than logloss (Ferri et al., 2009), hence our choice
 543 of logloss as a performance metric to compare different models coming from different training datasets.

544 3 RESULTS

545 3.1 Spatio-temporal ensemble machine learning framework

546 Table 4 shows the logloss performances based on the spatial 5-fold cross-validation repeated 5 times. This
 547 clearly indicates that GLM and Random forest have the highest predictive performances for all species.
 548 Both algorithms scored the lowest logloss among the other algorithms in 7 cases out of 14, with GLM
 549 excelling in potential distribution tasks and Random forest having better performances in the realized
 550 distribution tasks.

Table 4. Average logloss and standard deviation for the compared algorithms and for the subset of seven target species. In bold are highlighted the best performing learners for each task.

Distribution	Species	ANN	C5.0	GBT	GLM	KNN	RF	CART
Potential	<i>Abies alba</i>	0.170±0.014	0.042±0.003	0.089±0.001	0.009±0.001	0.057±0.012	0.023±0.001	0.057±0.005
Potential	<i>Castanea sativa</i>	0.253±0.018	0.055±0.004	0.100±0.003	0.015±0.002	0.066±0.016	0.028±0.003	0.082±0.011
Potential	<i>Fagus sylvatica</i>	0.423±0.014	0.069±0.005	0.115±0.002	0.021±0.001	0.076±0.009	0.036±0.002	0.170±0.008
Potential	<i>Picea abies</i>	0.450±0.012	0.070±0.002	0.120±0.001	0.032±0.001	0.091±0.008	0.040±0.002	0.142±0.008
Potential	<i>Pinus halepensis</i>	0.341±0.033	0.035±0.001	0.087±0.002	0.008±0.001	0.049±0.010	0.016±0.003	0.076±0.017
Potential	<i>Pinus sylvestris</i>	0.451±0.016	0.080±0.004	0.134±0.001	0.050±0.001	0.107±0.005	0.043±0.003	0.212±0.006
Potential	<i>Quercus robur</i>	0.401±0.018	0.074±0.006	0.120±0.002	0.023±0.001	0.070±0.009	0.040±0.003	0.156±0.006
Realized	<i>Abies alba</i>	0.628±0.004	0.143±0.002	0.138±0.007	0.074±0.007	0.206±0.039	0.069±0.007	0.189±0.016
Realized	<i>Castanea sativa</i>	0.565±0.019	0.247±0.074	0.187±0.020	0.161±0.018	0.372±0.068	0.107±0.021	0.271±0.035
Realized	<i>Fagus sylvatica</i>	0.544±0.013	0.094±0.017	0.113±0.003	0.050±0.005	0.145±0.034	0.048±0.006	0.127±0.016
Realized	<i>Picea abies</i>	0.650±0.058	0.122±0.009	0.168±0.003	0.093±0.001	0.244±0.022	0.079±0.005	0.244±0.015
Realized	<i>Pinus halepensis</i>	0.595±0.028	0.125±0.034	0.127±0.004	0.039±0.004	0.112±0.033	0.053±0.009	0.143±0.027
Realized	<i>Pinus sylvestris</i>	0.316±0.016	0.072±0.005	0.114±0.003	0.056±0.006	0.103±0.017	0.042±0.006	0.111±0.011
Realized	<i>Quercus robur</i>	0.627±0.039	0.183±0.024	0.180±0.008	0.107±0.009	0.249±0.032	0.098±0.007	0.248±0.024

551 Overall, GLM shows best performance in sense of achieving the lowest values of standard deviation.
 552 Random forest, however, scored close to GLM hence the difference does not seem to be significant.
 553 The absolute difference between values scored by GLM and Random forest is lower (0.089) than when
 554 Random forest had the advantage over GLM (0.105). This indicates a high reliability of Random forest
 555 performances even when other models outperform it. The neural network scored the highest logloss values
 556 in all tasks, so it was immediately excluded from the pool of level 0 models to choose from. It was time
 557 consuming to find a common hyperparameter range well suited for different tasks, since neural networks
 558 are often extremely situation-dependent. After a preliminary selection, we used the range shown in Table

3: despite that, our results remained inferior to those obtained with the other learners. On top of that, the *mlr* implementation of neural networks, based on the *deepnet* package, doesn't allow the use of ReLU (rectified linear activation function) as an activation function, which would have been beneficial for our purposes. Based on logloss performances, we selected RF and GLM as the first two components of the ensemble. Based on similar values of logloss (within one standard deviation of the average performance) scored by C5.0, GBT, KNN and CART, we used computational costs to choose the third component model. Table 5 shows the required time by each algorithm to run the hyperparameter tuning, 5-fold spatial cross validation and to predict a probability map on one tile. KNN was excluded due to computing time values being from one to two order of magnitude higher than the ones scored by the other models.

Table 5. Hyperparameter tuning, cross validation and prediction time for each model and distribution task. Time values are reported in seconds. Tests were conducted in a parallel computing setup on a CPU server running 2 x Intel(R) Xeon(R) Gold 6248R - 3.00GHz (96 threads) with 504 GB RAM.

Distribution	Process	ANN	C5.0	GBT	GLM	KNN	RF	Rpart
Potential	Tuning	661.2	220.7	680.1	345.9	2334.6	318.3	1019.3
Potential	Cross validation	54.8	45.5	34.4	305.7	582.5	426.1	14.5
Potential	Prediction	24.12	231.1	24.8	14.8	19272.9	53.4	14.1
Realized	Tuning	1772.3	944.5	851.1	427.2	16198.4	1873.8	3298.1
Realized	Cross validation	16.2	184.8	20.5	344.1	2919.9	2347.5	39.9
Realized	Prediction	19.3	389.3	32.6	17.8	> 1 day	78.4	16.4
Total		2547.92	2015.9	1643.5	1455.5	> 1 day	5097.5	4402.3

Even though the CART scores very low value in cross validation and prediction time in both potential and realized tasks, tuning time is the second highest, just behind KNN. C5.0 is faster than GBT in the whole potential workflow (497.3 seconds against 739.3) but slower in the realized workflow (1518.6 seconds against 904.2). Considering both workflows, GBT proved to be faster and more consistent in cross validation and prediction time, showing an increase in tuning time of just 25% with double the amount of training data (see Table 2). After the training of the component models and the meta-learner, this is an example of final fitted ensemble model:

```
Species: Pinus_sylvestris
Distribution: Potential
Logloss: 0.0255

Call:
stats::glm(formula = f, family = "binomial", data = getTaskData(.task,
  .subset), weights = .weights, model = FALSE)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-3.5683  0.0586  0.0587  0.0595  3.4006

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)   -5.78660    0.05221  -110.83  <2e-16 ***
classif.ranger  8.12492    0.10953   74.18  <2e-16 ***
```

```

classif.xgboost  1.49208    0.08233    18.12    <2e-16 ***
classif.glmnet  2.53462    0.07296    34.74    <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 290519  on 459919  degrees of freedom
Residual deviance: 23472  on 459916  degrees of freedom
AIC: 23480

Number of Fisher Scoring iterations: 8

```

575 The output shows that all algorithms are significant predictors for the response variable (column $Pr(>|z|)$),
576 while from the coefficient estimate values we infer that predictions coming from Random forest have, for
577 this specific case, the highest weight in the final predictions of the meta-learner, followed by GLM with
578 Lasso and finally GBT.

579 3.2 Variable importance

580 We generated the variable importance for the three component models for all species and distributions:
581 for each of these combinations we computed the relative variable importance and then took the top-20
582 variables from each set. This resulted in a total of 146 unique variables for potential distributions and 178
583 for realized distributions; of those, 54 variables are present in both potential and realized set of important
584 variables. To better analyze the results and identify driving factors of species distributions, we aggregated
585 the whole set of 270 variables in 7 macro-classes:

- 586 • Climate (i.e. precipitation, wind speed, water vapor, snow probability etc.),
- 587 • Temperature (i.e. time series of recorded temperatures for the observed time period),
- 588 • Bioclim (i.e. bioclimatic variables from CHELSA),
- 589 • Topography (i.e. DTM and DTM-derivative variables),
- 590 • Landsat band (i.e. all percentiles, all seasons),
- 591 • Distribution (i.e. species distribution maps from European Atlas of Forest Tree Species),
- 592 • Spectral index (i.e. spectral indices derived from Landsat bands).

593 Results are presented in Fig. 6. Each plot can be divided in four quadrants, from the top left clockwise:
594 variables with high relative importance but low frequency (i.e. important for one or few species), variables
595 with high importance and high frequency (i.e. important for all species), variables with low importance
596 and high frequency (i.e. they occurred often but were not important) and variables with low importance
597 and low frequency.

598 For potential distribution, diffuse irradiation, precipitation of the driest quarter (BIO17) and precipita-
599 tion of the driest month (BIO14) were the most important and most frequent predictors across all models

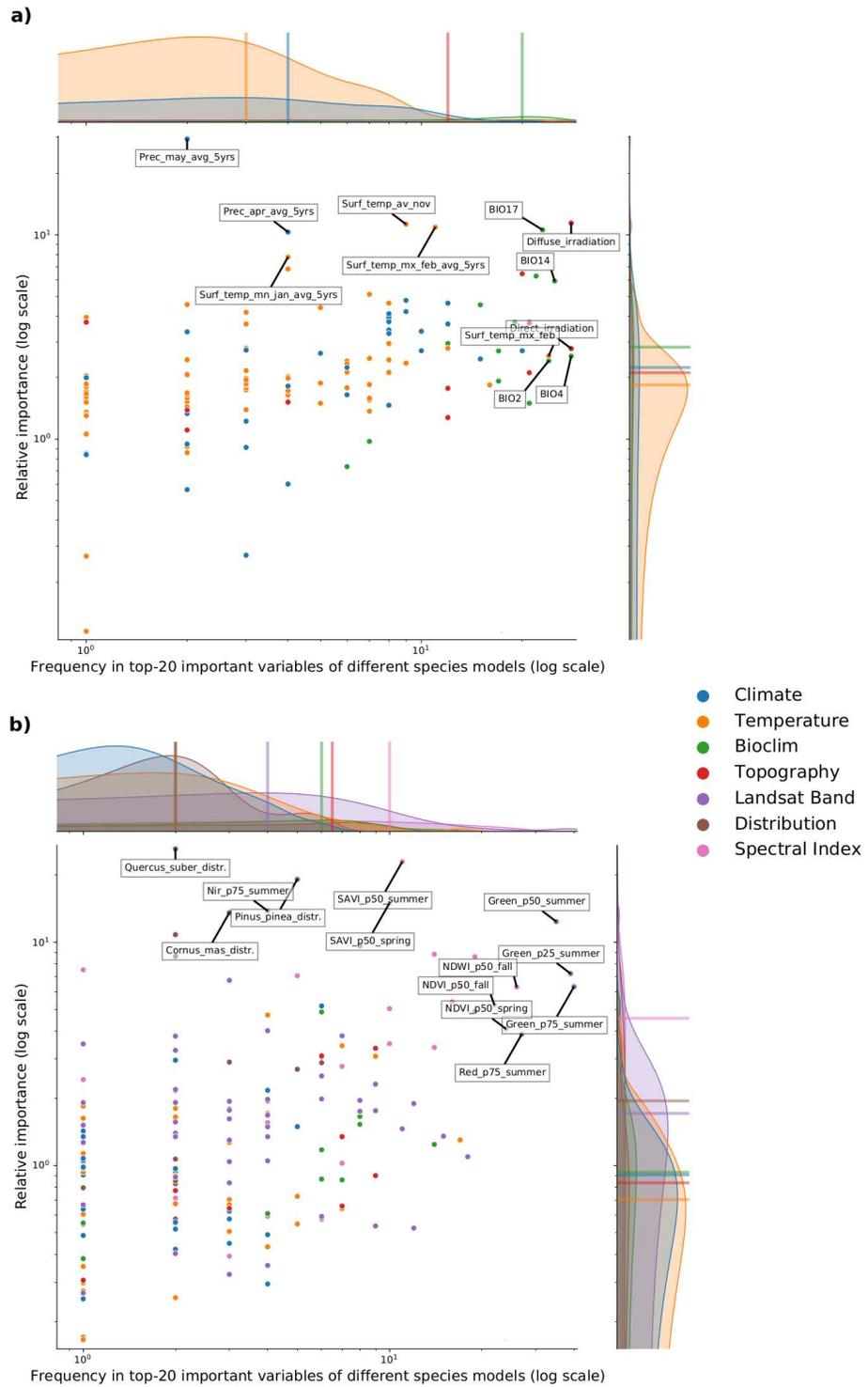


Figure 6. Relative variable importance vs frequency of the variables of the top-20 most important across the component models and all species for potential **(a)** and realized **(b)** distribution. Variables were aggregated in seven different classes, colors are shown in legend. Labeled dots are variables that recorded high values of relative variable importance or frequency: each plot shows in the top right corner the most important and most frequent variables across all models and species.

600 and species. The density distributions per macro-class help understanding how the Bioclim macro-class
601 was the one with on average both most important and most frequent variables. Other variables are more
602 species-specific: the cumulative precipitation of May computed on a 5 years time window records the
603 highest absolute value in relative importance but it was important for just two species (*Abies alba* and
604 *Castanea sativa*, see additional statistical outputs at <https://zenodo.org/record/5821865>). The
605 Temperature macro-class accounts the highest numbers of predictors, but the values recorded in both
606 variable importance and frequency are the lowest among all the macro-classes. The Climate macro-class
607 had the largest variety in predictors and variables in this class are homogeneously spread out across all the
608 species in both variable importance and frequency.

609 For realized distribution, the summer aggregates of Landsat green (25th and 50th quantiles) were
610 the two most important and most frequent variables across all models and species, closely followed
611 by the fall aggregates of NDVI and NDWI and summer and spring aggregates of SAVI. Reflectance-
612 based macro-classes clearly outperformed the other ones in this case. The distribution maps scored
613 the highest values for variable importance (distribution of the *Quercus suber* and the *Pinus pinea*) but
614 they were species-specific (i.e. important for just three species, see additional statistical outputs at
615 <https://zenodo.org/record/5821865>). Despite the green Landsat band scoring the best values
616 across all the other predictors, the Spectral index macro-class scored on average as the most important in
617 both variable importance and frequency.

618 Overall, the component models show more differences in variable importance in the potential distribu-
619 tion models than in the realized ones. On average, Random forest and Gradient Boosted Trees selected the
620 same variables in the top-10 but not always in the same order, while GLM tended to choose completely
621 different variables. This suggests how the ensemble models tend to use a wider proportion of the feature
622 space than single models. This tendency is most apparent in the potential distribution models. In the
623 realized distribution models, the component models agree in selecting the top-10 most important variables
624 predictors from Landsat bands or Spectral indices. Random forest and Gradient Boosted Trees considered
625 on average the Landsat bands as the most important, while GLM selected the spectral indices more often.

626 3.3 Accuracy assessment

627 We ran a 5-fold spatial cross validation repeated 5 times with the tuned component models and the
628 ensemble model to provide conservative estimates of predictive performance. Predictions on the left-out
629 folds were aggregated to assess model performance. Table 6 shows logloss and standard deviation values
630 for component models and ensemble model across all species and distributions compared with the value
631 of logloss for the intercept-only model. All models were better than random assignment. In general,
632 the models for realized distribution achieved better predictive performances than those for potential
633 distribution. This does not hold for three species: *Fagus sylvatica*, *Olea europaea* and *Quercus suber*.
634 The ensemble model:

- 635 • **outperformed** the component models in **15** cases of the potential distribution tasks but only in **one**
636 case of the realized distribution tasks,
- 637 • performed **as good as** the best component model in **one** case of the potential distribution tasks and

638 in **five** cases of the realized distribution tasks,

- 639 • performed **worse** than the best component model in **ten** cases of the realized distribution tasks,
- 640 with small differences.

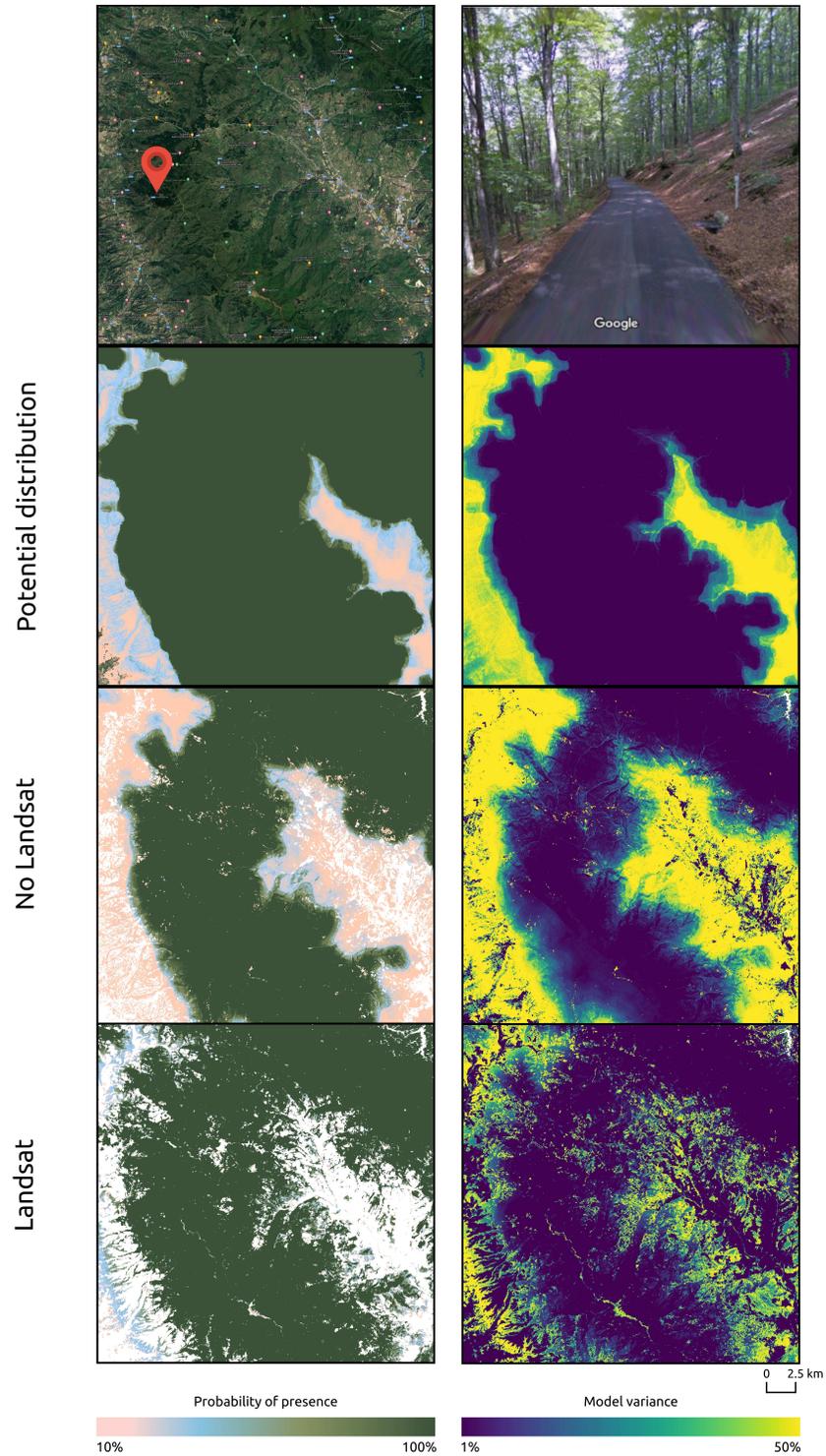
Table 6. Average logloss and standard deviation for the component learners and the ensemble model. In bold are highlighted the best performing learners for each task. Random logloss values are dependent on presence-absence ratio in the species dataset and are here used as a baseline for predictive performances comparison.

Distribution	Species	GBT	GLM	RF	Ensemble	Random
Potential	<i>Abies alba</i>	0.062 ±0.008	0.071 ±0.005	0.071 ±0.011	0.046 ±0.005	0.606
Potential	<i>Castanea sativa</i>	0.133 ±0.052	0.148 ±0.016	0.119 ±0.027	0.080 ±0.013	0.622
Potential	<i>Corylus avellana</i>	0.109 ±0.018	0.125 ±0.008	0.118 ±0.027	0.080 ±0.009	0.638
Potential	<i>Fagus sylvatica</i>	0.032 ±0.004	0.074 ±0.007	0.048 ±0.008	0.029 ±0.003	0.544
Potential	<i>Olea europaea</i>	0.007 ±0.004	0.005 ±0.002	0.008 ±0.004	0.005 ±0.002	0.233
Potential	<i>Picea abies</i>	0.069 ±0.004	0.089 ±0.005	0.072 ±0.005	0.067 ±0.004	0.634
Potential	<i>Pinus halepensis</i>	0.048 ±0.012	0.037 ±0.005	0.057 ±0.011	0.033 ±0.006	0.596
Potential	<i>Pinus nigra</i>	0.091 ±0.011	0.101 ±0.008	0.117 ±0.012	0.077 ±0.005	0.494
Potential	<i>Pinus pinea</i>	0.033 ±0.007	0.035 ±0.004	0.047 ±0.012	0.023 ±0.004	0.434
Potential	<i>Pinus sylvestris</i>	0.031 ±0.003	0.055 ±0.003	0.035 ±0.004	0.026 ±0.002	0.315
Potential	<i>Prunus avium</i>	0.096 ±0.010	0.115 ±0.007	0.103 ±0.019	0.071 ±0.006	0.509
Potential	<i>Quercus cerris</i>	0.063 ±0.017	0.037 ±0.009	0.049 ±0.011	0.024 ±0.005	0.360
Potential	<i>Quercus ilex</i>	0.071 ±0.008	0.094 ±0.007	0.071 ±0.015	0.058 ±0.007	0.508
Potential	<i>Quercus robur</i>	0.075 ±0.007	0.113 ±0.010	0.084 ±0.006	0.069 ±0.007	0.626
Potential	<i>Quercus suber</i>	0.016 ±0.012	0.011 ±0.004	0.021 ±0.012	0.009 ±0.003	0.348
Potential	<i>Salix caprea</i>	0.106 ±0.012	0.138 ±0.010	0.117 ±0.011	0.091 ±0.008	0.662
Realized	<i>Abies alba</i>	0.017 ±0.002	0.014 ±0.002	0.022 ±0.002	0.016 ±0.002	0.264
Realized	<i>Castanea sativa</i>	0.024 ±0.004	0.026 ±0.003	0.031 ±0.004	0.025 ±0.002	0.281
Realized	<i>Corylus avellana</i>	0.024 ±0.003	0.02 ±0.002	0.031 ±0.002	0.022 ±0.002	0.212
Realized	<i>Fagus sylvatica</i>	0.030 ±0.002	0.031 ±0.001	0.043 ±0.001	0.030 ±0.002	0.554
Realized	<i>Olea europaea</i>	0.012 ±0.004	0.006 ±0.001	0.012 ±0.004	0.008 ±0.002	0.107
Realized	<i>Picea abies</i>	0.037 ±0.002	0.043 ±0.002	0.049 ±0.001	0.040 ±0.002	0.666
Realized	<i>Pinus halepensis</i>	0.013 ±0.002	0.011 ±0.001	0.020 ±0.003	0.011 ±0.001	0.346
Realized	<i>Pinus nigra</i>	0.018 ±0.002	0.016 ±0.002	0.027 ±0.004	0.016 ±0.002	0.297
Realized	<i>Pinus pinea</i>	0.013 ±0.005	0.009 ±0.002	0.019 ±0.005	0.010 ±0.002	0.245
Realized	<i>Pinus sylvestris</i>	0.043 ±0.001	0.053 ±0.001	0.058 ±0.001	0.045 ±0.002	0.682
Realized	<i>Prunus avium</i>	0.022 ±0.002	0.020 ±0.002	0.030 ±0.002	0.020 ±0.002	0.165
Realized	<i>Quercus cerris</i>	0.009 ±0.002	0.006 ±0.001	0.013 ±0.003	0.007 ±0.002	0.113
Realized	<i>Quercus ilex</i>	0.023 ±0.004	0.019 ±0.002	0.027 ±0.004	0.018 ±0.002	0.291
Realized	<i>Quercus robur</i>	0.041 ±0.002	0.036 ±0.002	0.045 ±0.003	0.036 ±0.001	0.448
Realized	<i>Quercus suber</i>	0.012 ±0.008	0.005 ±0.001	0.011 ±0.004	0.006 ±0.002	0.180
Realized	<i>Salix caprea</i>	0.028 ±0.002	0.031 ±0.002	0.035 ±0.002	0.030 ±0.002	0.268

641 Among the best performing models, the ensemble for potential distribution of *Olea europaea* and GLM
642 for realized distribution of *Quercus suber* achieved both the absolute and relative (i.e. scaled with the
643 random value) lowest values of logloss. The ensemble for potential distribution of *Salix caprea* achieved
644 the highest absolute value, while the highest relative value was achieved by the ensemble for potential
645 distribution of *Pinus nigra*.

646 3.4 Influence of high resolution on predictive performances

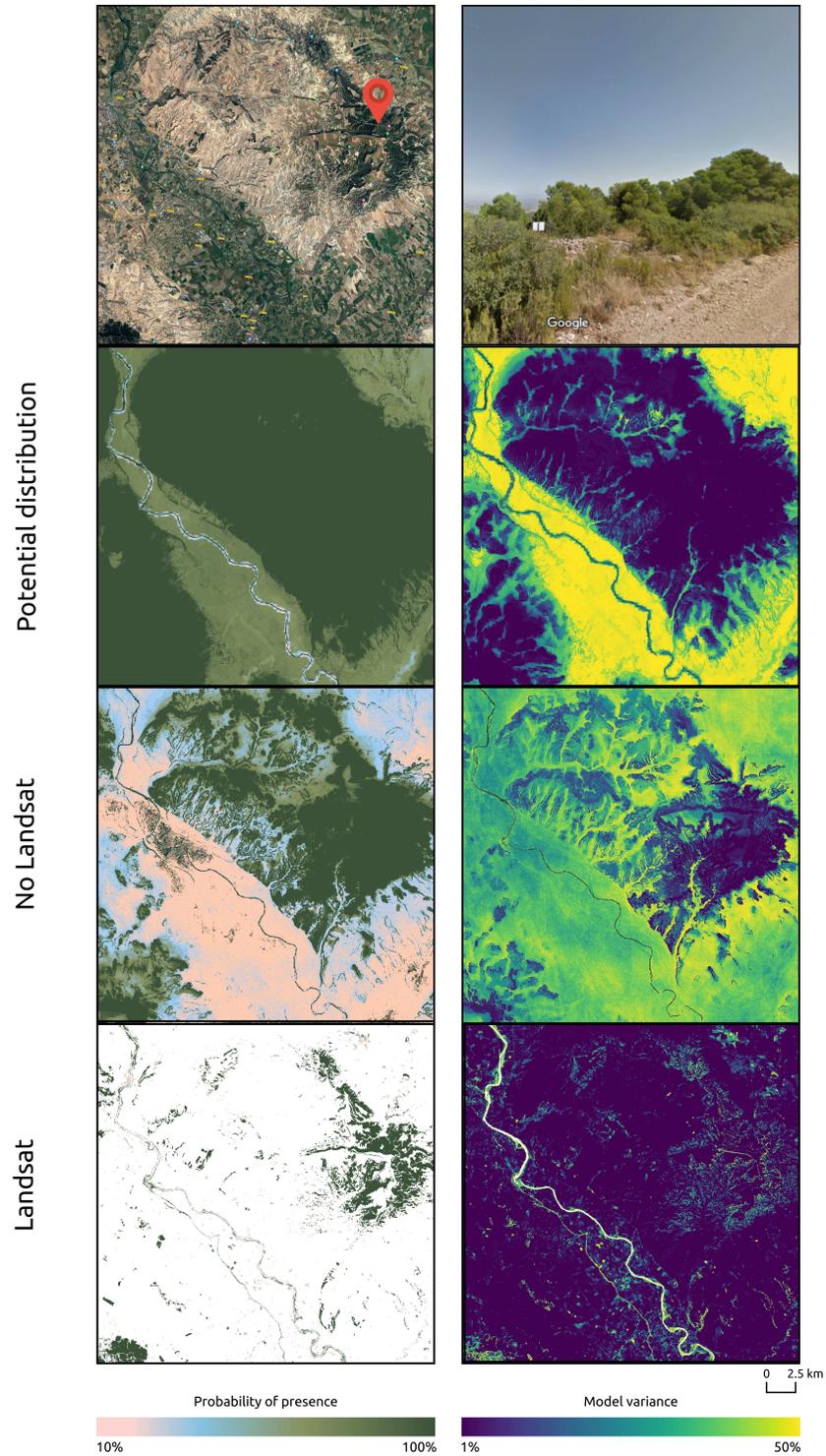
647 To assess the effect of high resolution products on predictive performances, we excluded Landsat bands
648 and Landsat-derived spectral indices from the list of predictors used for realized distribution. We then



649

Figure 7. Comparison of potential distribution with realized distribution for the period 2018–2020 for *Fagus sylvatica* on a 30 km tile located in Tuscany, Italy. This area was chosen due to the presence of a national park (Foreste Casentinesi National Park) and a Natura 2000 area (Vallombrosa and Sant’Antonio forest) well known for *Fagus sylvatica* stands. Satellite view provided by Google Maps (left) and ground view provided by Google Street View (right) are shown in Row 1. Row 2 to row 4 show probability (left) and uncertainty (right) maps. Two realized distributions are shown: excluding and including Landsat data among the predictor variables.

650



651

652

Figure 8. Comparison of potential distribution with realized distribution for the period 2018–2020 for *Pinus halepensis* on a 30 km tile located in the Ebro Basin, Spain. The species was introduced in this area in the Mesolithic period. Satellite view provided by Google Maps (left) and ground view provided by Google Street View (right) are shown in Row 1. Row 2 to row 4 show probability (left) and uncertainty (right) maps. Two realized distributions are shown: excluding and including Landsat data among the predictor variables

653 applied our spatio-temporal ensemble machine learning framework (feature selection — hyperparameter
654 tuning — ensemble model training) on each species and we ran a 5-fold spatial cross validation repeated
655 5 times to evaluate model performances. For the ensemble model we used the same component models
656 (Random forest, GBT and penalized GLM) and meta-learner (penalized logistic regression). Results
657 of this analysis were then compared with the performances achieved by the ensemble models using
658 Landsat data. Comparison is shown in Table 7: for all species, ensemble models including Landsat data
659 consistently overperformed models of the same species without the Landsat data. In all cases, models
660 including Landsat data show better predictive performances (lower values of logloss) and less uncertainty
661 (lower values of standard deviation). Fig. 7 and 8 show a visual comparison of the models for two species
662 (*Fagus sylvatica* and *Pinus halepensis*) on sample areas.

Table 7. Average logloss for modelling realized distribution with and without the Landsat bands and spectral indices. Random logloss values are shown as a baseline for predictive performances.

Species	Landsat	No Landsat	Random
<i>Abies alba</i>	0.016 ±0.002	0.036 ±0.003	0.264
<i>Castanea sativa</i>	0.025 ±0.002	0.065 ±0.007	0.281
<i>Corylus avellana</i>	0.022 ±0.002	0.045 ±0.003	0.212
<i>Fagus sylvatica</i>	0.030 ±0.002	0.088 ±0.003	0.554
<i>Olea europaea</i>	0.008 ±0.002	0.010 ±0.003	0.107
<i>Picea abies</i>	0.040 ±0.002	0.093 ±0.002	0.666
<i>Pinus halepensis</i>	0.011 ±0.001	0.041 ±0.004	0.346
<i>Pinus nigra</i>	0.016 ±0.002	0.039 ±0.004	0.297
<i>Pinus pinea</i>	0.010 ±0.002	0.018 ±0.003	0.245
<i>Pinus sylvestris</i>	0.045 ±0.002	0.104 ±0.003	0.682
<i>Prunus avium</i>	0.020 ±0.002	0.041 ±0.003	0.165
<i>Quercus cerris</i>	0.007 ±0.002	0.018 ±0.003	0.113
<i>Quercus ilex</i>	0.018 ±0.002	0.032 ±0.003	0.291
<i>Quercus robur</i>	0.036 ±0.001	0.108 ±0.006	0.448
<i>Quercus suber</i>	0.006 ±0.002	0.010 ±0.003	0.180
<i>Salix caprea</i>	0.030 ±0.002	0.051 ±0.003	0.268

663 4 DISCUSSION

664 4.1 Summary findings

665 In this paper we have developed, tested and reported a methodological framework for predicting the
666 potential and realized distributions of 16 forest tree species by Ensemble Machine Learning and analysis-
667 ready Earth Observation data. In general, our ensemble model achieved better predictive performances
668 than individual models when modeling the potential distribution, while performing slightly worse than
669 the best component model for ten species when modeling the realized distribution.

670 Even at high resolution (i.e. on a local scale), climate proved to be the key driver of vegetation
671 potential distribution across Europe mainly through temperature and precipitation. These results are
672 consistent with findings of Hutchinson (1957) and Mather and Yoshioka (1968). Reflectance-based
673 covariates were the most important predictors of the realized distributions. Overall, potential distribution
674 proved to be significantly more complex to map accurately than realized distribution. In general, the
675 ensemble and component models achieved better predictive performances for the realized distributions

676 than for the potential distributions as judged from the cross-validation logloss estimates (see Table 6 and
677 Fig. 9).

678 Our results indicate a consistent increase in predictive performances for realized distribution when
679 adding high resolution data, especially Landsat data at 30 m resolution and vegetation indices to the list
680 of predictors (see Table 7). Significant findings of our work include:

- 681 1. Predictive modeling of forest tree species can be efficiently automated to the level of full automation,
682 but this assumes high quality / artifact free training points with a homogenous distribution of
683 occurrence and absence points whenever possible (Fig. 3)
- 684 2. Complexity of Ensemble Machine Learning methods can be significantly reduced by implementing
685 the following two generic processes: (1) hyperparameter tuning on data subsets, (2) running efficient
686 feature selection (Fig. 4).
- 687 3. Analysis-ready Landsat time-series are maybe cumbersome to prepare and gap-fill for clouds and
688 artifacts (Witjes et al., 2021), but overall come as the most important inputs for mapping forest tree
689 species, hence the Landsat project stands out as the most important data source for continuous and
690 consistent time-series forest mapping.

691 Compared to the previous application of Machine Learning methods for mapping forest tree species
692 over the whole European continent described in Hengl et al. (2018), we have identified the following
693 methodological improvements:

- 694 • A framework to face the problem of preferential sampling and absence data unavailability to create a
695 presence-absence datasets mainly by using LUCAS points (EUROSTAT, 2017). In contrast, Hengl
696 et al. (2018) and de Rigo et al. (2016a) used only occurrence points, which probably introduced
697 bias (i.e. overestimation or underestimations) in the predictions.
- 698 • A framework for spatio-temporal predictive mapping. We built a single model per forest tree
699 species to predict through the whole spatio-temporal datacube cube and produce time-series of
700 predictions. Maps can then be used as input data for further analysis / work, i.e. to detect species
701 geographical range contraction or expansion, areas of forest degradation / restoration, species
702 composition changes.
- 703 • Methodological steps to help with reduction of the model complexity (Fig. 4): these have shown to
704 be beneficial for reducing the production costs and enabled us to map more forest tree species.

705 4.2 Modeling framework

706 Combining models using the ensemble approach is thought to reduce model uncertainty and increase its
707 robustness in modelling species distributions (Araújo and New, 2007). We used ensemble with stacked
708 generalization as ensemble approach, which has not been tested yet for species distribution modelling.
709 We also trained the models in a spatio-temporal framework, expecting the models to generalize better
710 when predicting in a temporal window not included in the training data. Part of the intent of the paper
711 is to provide a reproducible framework to model species distributions, so we compare our results with

712 those available from previous publications that use ensemble modeling for either potential or realized
713 distributions. Our study, however:

- 714 • models both potential and realized distributions,
- 715 • uses only one functional group (trees) on one continent (Europe),
- 716 • adds several layer of complexity (ensemble based on stacking, spatio-temporal framework),
- 717 • uses a large (585) set of predictor variables,
- 718 • uses logloss as a performance measure,
- 719 • uses spatial cross validation,
- 720 • does not test model transferability,

721 so not all results from other studies can be directly compared. Hao et al. (2020) used a similar method-
722 ological framework to the one used in this study. They modelled the distribution of 13 species of the genus
723 *Eucalyptus* in South Australia and tested performances of ensemble model against individual models;
724 they used mean and weighted average as ensemble strategies. They also tested cross validation versus
725 spatial cross validation for model performances. The study doesn't specify which type of distribution was
726 modelled: according to the definition provided in this study and the list of predictor variables used by Hao
727 et al. (2020), we can classify the task as potential distribution.

728 Their results show how spatial cross validation performances were more conservative than cross
729 validation ones when compared with performances on independent validation sets. This supports and
730 reinforces our use of spatial cross validation as a validation strategy for the modeling framework. Ensemble
731 models performed well but were outperformed by not tuned individual models and by a tuned GBT. There
732 was also no clear advantage in predictive performances when using different ensemble strategies. This is
733 in contrast with our results, where the ensemble based on stacking outperformed even tuned component
734 models in 15 cases of the potential distribution tasks and performing as good as the best component model
735 in the last case. This may be an indication of stacking being a better ensemble strategy when modeling
736 species distribution.

737 In our case, a tuned GBT outperformed the ensemble only in realized distribution tasks. By comparing
738 Table 4 and 6 we can also see that Random forest and GLM are clearly the best component models to
739 map both potential and realized distributions when trained on a data sample, but GBT often outperforms
740 Random forest or even the ensemble when tuned and trained on the whole dataset. In general, differences
741 in predictive performances between the ensemble and the component models are also higher in potential
742 distribution than in realized distribution. The list of variable importance per component model, species
743 and task (see additional statistical outputs at <https://zenodo.org/record/5821865>) may give an
744 insight to this: in the potential tasks, the component models use different parts of the feature space before
745 the predictions are combined by the meta-learner. All the models select as most important variables for the
746 task different predictors. For realized distribution tasks, the models all agree in selecting either Landsat

747 bands or spectral indices as most important variables, resulting in predictions that are highly correlated
748 and with less variance between the models.

749 Ensemble modeling is known to perform best when there is a high diversity between the base
750 models (a) and no or negative correlation between their outputs (b) (Zhou, 2019). The introduction
751 of Landsat bands and spectral indices in general greatly increased the predictive performances of the
752 models for realized distribution (see Table 7) compared to potential distribution models (see Table
753 6). However, it also homogenized predictions, which makes the (b) condition reported above not
754 valid. We separately compared the repeated spatial cross validation performances of ensemble and
755 component models excluding the Landsat bands and spectral indices (see additional statistical outputs at
756 <https://zenodo.org/record/5821865>). In this case, the ensemble performed worse than the best
757 component model only in four cases instead of ten. While most of ensemble SDM studies use no more
758 than 15–20 predictor variables (Kaky et al., 2020; Hao et al., 2019; Kindt, 2018), even after feature
759 selection we ended up using more than 300 variables to model the realized distribution of some species.
760 A stricter feature selection could maybe remove the problem in future works. In general, if the ensemble
761 provides predictive performances as good as or worse than the best component model, the best component
762 model must be preferred (Zhang and Ma, 2012). However, ensemble models can still provide more
763 advantages than individual models since they reduce model uncertainty and are more robust towards
764 extrapolation (Mehra et al., 2019).

765 4.3 Species distributions

766 Our cross-validation accuracy assessment results indicate high predictive performances for all species, in
767 both potential and realized distributions. In the case of mapping potential distribution, diffuse irradiation
768 and precipitation of the driest quarter (BIO17) come as overall most important covariate layers. By design,
769 potential distribution maps produced in this study are indicators of where a species is not likely to survive
770 (see section 2.1). Our interest was to map those areas deemed not suitable by the model and identify
771 which of the considered predictor variables would have been the limiting factors. Which environmental
772 variables and their relative importance as limiting factors are still unclear for many tree species and few
773 peer reviewed studies focused on investigating the importance of different environmental predictors on
774 potential distributions are available. We found that, on average, each component model considers two
775 or more predictors from the Bioclim macro-class among the top-10 most important variables to predict
776 the potential distribution (see additional statistical outputs at <https://zenodo.org/record/5821865>).
777 Previous findings in literature have shown the importance of bioclimatic variables when modeling species
778 distributions (Fourcade et al., 2018), but this may also be a consequence of bioclimatic variables and
779 elevation being the most employed, if not the only, predictors in numerous SDM studies (Fois et al.,
780 2018). Bucklin et al. (2015) compared the influence of different sets of environmental predictors on model
781 performances, but the list of predictors used in the study included human influenced factors, so their
782 results cannot be used to assess the driving factors for potential distributions. Pearson and Dawson (2003)
783 argued that environmental predictors only are not able to describe all the processes that limit species
784 ranges, but this claim is directed at realized distributions only. Even if our results show the bioclimatic
785 variables as the most important predictors for potential distributions, further studies in this direction may

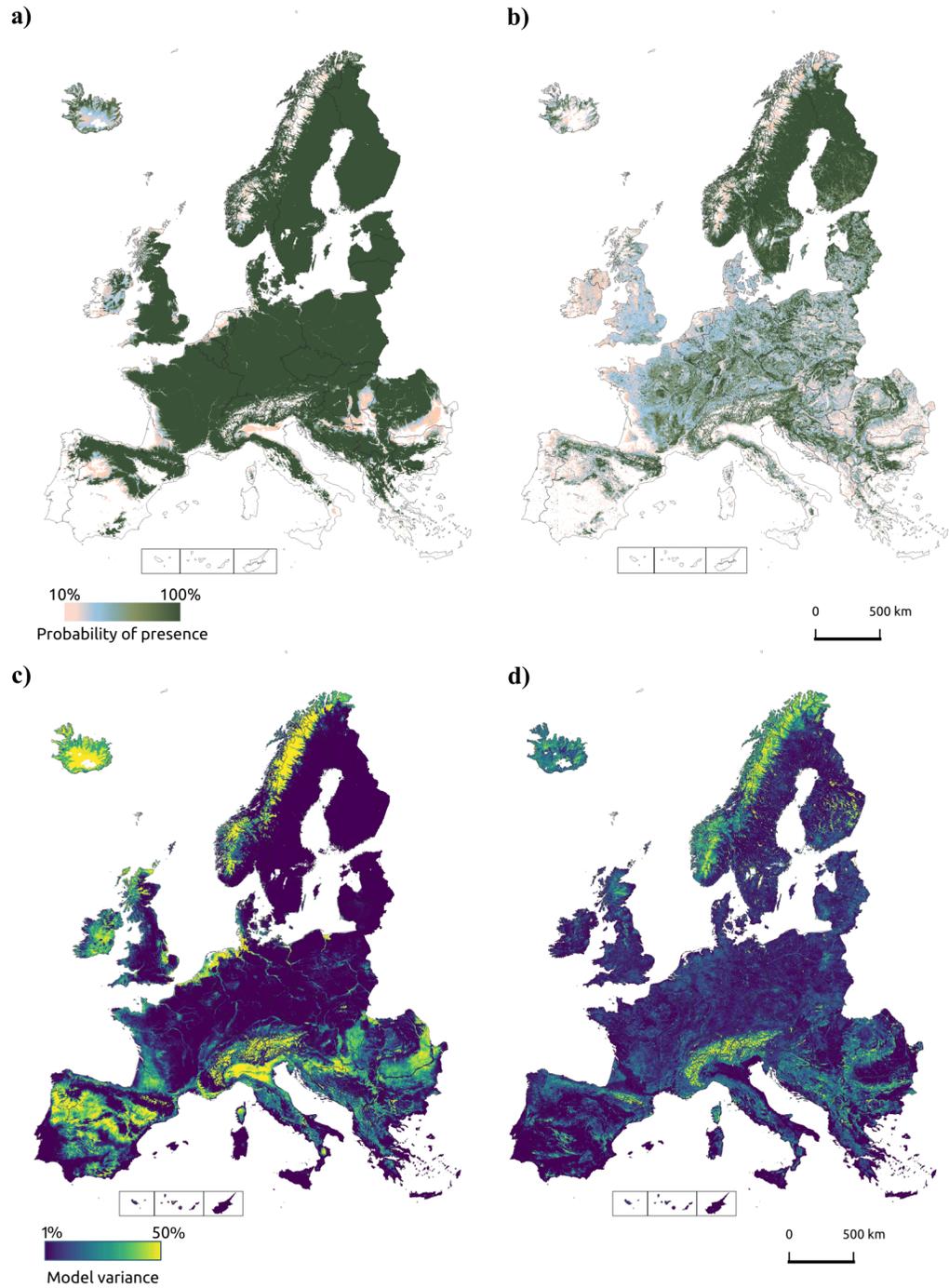


Figure 9. Potential and realized distribution for *Pinus sylvestris* for the period 2018–2020. Clockwise, from the top left corner: probability map for potential distribution (a), probability map for realized distribution (b), uncertainty map for potential distribution (c) and uncertainty map for realized distribution (d).

786 be needed. The scale of the study may affect the importance of predictor variables: on a large scale,
787 distribution may be influenced by macro environmental factors, while at a local scale, other environmental
788 factors may limit mangrove distribution more significantly. [Walther and Meier \(2017\)](#) and [Weigel et al.](#)
789 [\(2019\)](#) proved that soil properties are more important than either bioclimatic or only climatic variables
790 when modeling potential tree species distribution at, respectively, country and regional scale.

791 For realized distribution, variable importance confirms that Earth Observation layers such as the 25th
792 and 50th quantile summer aggregates for the Landsat green band and the 50th quantile fall aggregates of
793 NDVI and NDWI are overall the most important layers for mapping realized distribution of species (Fig.
794 6). The inclusion of Landsat data and derived spectral indices increases predictive performances (see Table
795 7) and contains more detailed information on species distribution ranges (see Fig. 7 and 8). Importance of
796 NDVI is well known since it is one of the most used proxies in vegetation studies such as biodiversity
797 estimation ([Madonsela et al., 2017](#); [He et al., 2009](#)), net primary productivity ([Schloss et al., 1999](#)) and
798 land degradation ([Easdale et al., 2018](#)), phenology ([Fawcett et al., 2021](#)) and species composition changes
799 ([Wang et al., 2021](#)). NDVI incorporates information from the red and the near-infrared (NIR) portion of
800 the electromagnetic spectrum (see Table 1). Vegetation's behavior in this portion of the spectrum has long
801 been used in vegetation mapping to distinguish between coniferous and deciduous tree species ([Hoffer,](#)
802 [1984](#)). The green band, although usually less important than the red and NIR band, has already proved
803 useful in vegetation mapping to classify forest types ([Gao et al., 2015](#)), predict forest variables (stem
804 volume, diameter and tree height) at species level ([Astola et al., 2019](#)) and forest biomass at community
805 level ([Nandy et al., 2017](#)).

806 Comparing our results with chorological maps from the European Atlas of Forest Tree Species
807 ([San-Miguel-Ayanz, J., de Rigo, D., Caudullo, G., Houston Durrant, T., Mauri, 2016](#)), we can see that in
808 general both potential and realized distribution correctly capture the species ranges. Overall, potential
809 distribution maps show homogeneous patterns of high probability values for all target species, while
810 realized distribution maps show very fragmented patterns. Fig. 7 and 8 show sample areas for two
811 different species where the potential distribution model indicates high probability for the target species.
812 The realized distribution model helps discriminating the presence or absence of the species due to biotic
813 or other external factors. A high geographical overlap between probability maps of realized distribution
814 may reflect co-existence of different species within the same forest stands and could help with clearly
815 defining forest communities.

816 **4.4 High resolution contributions: is finer always better?**

817 Bioclimatic variables available only at coarse spatial resolution were used as predictor variables in
818 both potential and realized distribution. The Landsat bands and the spectral indices were not the only
819 high resolution layers used in this study: terrain and terrain-derived predictors were also included at
820 30 m resolution. However, regardless of the terrain data high resolution, the tree species potential
821 distribution patterns mostly reflect the original spatial resolution of the bioclimatic variables. Thus,
822 climate influences species distribution on European scale. Even though this might indicate that mapping
823 potential distributions at high resolution may not be necessary, it can still be useful for different case
824 studies. For example, comparing the difference, and hence mapping the gap, between potential and

825 realized distribution at the same fine scale, as we show in Fig. 7 and 8, may prove to be an invaluable tool
826 for both forest managers and conservation planners that work on the local level.

827 Potential distribution maps can be used to identify suitable areas for species in reforestation and
828 restoration programs; realized distribution maps can inform the forest managers on the presence or
829 absence of said species in those areas at a particular point in space and time (see Fig. 10). By removing
830 the biotic factors that limit the presence of the species in a potential reforestation site, using multiple
831 distribution maps and including expert knowledge on species synecology, structurally complex forest
832 stands could be planned and developed in a much more informed way. The opposite approach could be
833 used by conservation planners. Potential distribution is modelled by studying the relationship between a
834 species and the environmental conditions found in its native range, where the species is at equilibrium
835 (Jiménez-Valverde et al., 2011). Invasive species are usually more abundant and have greater performances
836 in the introduced range than in their native ranges (Hierro et al., 2005). This is due to the absence of
837 biotic factors that normally limit species distribution in their native range in the introduced range. Thus,
838 a species that occupies only 10% of its potential distribution in its native range may end up occupying
839 a bigger percentage of it in the introduced range. Estimation of potential distribution in the introduced
840 range that depends only on environmental factors are conservative by definition, potential distribution
841 maps may provide a good indication to conservation planners of how much the invasive species could
842 spread in the introduced range.

843 For realized distribution, including high resolution predictor variables in the model not only increases
844 predictive performances but also lowers overall (see Table 7) and local (see Fig. 7 and 8) values of
845 uncertainty. For forest management purposes, a large, consistent, standardized, long-term and high
846 resolution image collection such as the one provided by the Landsat program can help extending in space
847 and time information on tree species presence, composition and abundance. A spatial resolution of 30 m
848 is particularly well suited for national forest inventory applications: Strickland et al. (2020) derived
849 probability maps of forest tree species for a 25 years time period (1985–2010) using yearly Landsat
850 composites to extend missing information from the Canadian NFI and estimating changes in forest cover,
851 species composition and forest disturbances. The increasing availability of even higher-spatial resolution
852 satellite data from the European Copernicus program (i.e. Sentinel 1 and 2) and commercial providers (i.e.
853 Planet) can potentially further enhance predictions by including more data and a better spatial matching
854 of in-situ and satellite-derived information.

855 4.5 Technical limitations

856 The spatio-temporal Ensemble Machine Learning based on stacking is highly computational and requires
857 significant investment in High Performance Computing. The total computing time required to conduct
858 feature selection, hyperparameter tuning, model training, accuracy assessment and predictions of all
859 probability and uncertainty maps described for just 1 species exceeds 25,000 CPU hours based on the Intel
860 Xeon Gold chips, which within our infrastructure implied 19 days of non-stop computing with full capacity.
861 This estimation doesn't include training data and predictor variables pre-processing and harmonization.
862 We decreased processing costs by using Amazon S3 (Palankar et al., 2008) and full parallelization
863 of all processes including spatio-temporal overlay (see <http://eumap.readthedocs.org/>), hyper-

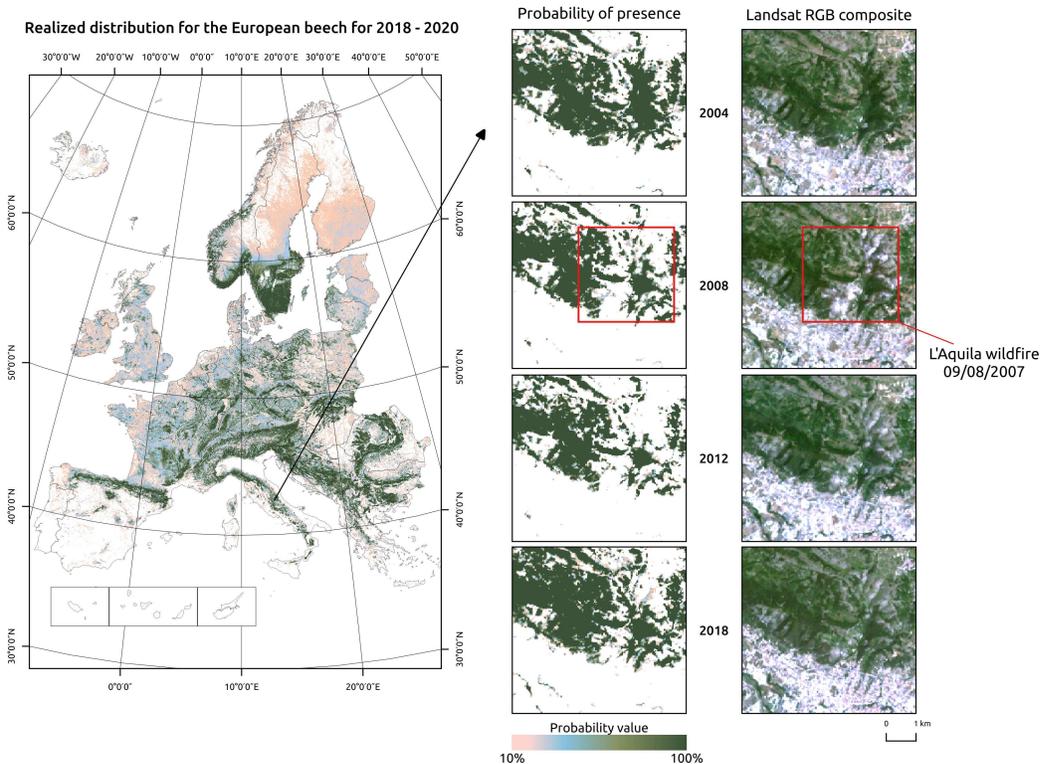


Figure 10. Realized distribution of *Fagus sylvatica* for the period 2018-2020. Detailed insets show a region around L'Aquila city, in Central Italy. The *Fagus sylvatica* forest on the northern outskirts of the city was affected by a serious wildfire in 2007. The realized distribution maps can be used to track changes through time.

864 parameter tuning and prediction, but the processing costs remain high. Due to the high production costs
 865 and data availability we have eventually decided to deviate from the original plan to map all forest tree
 866 species from the European Forest Atlas (76) to focus on a smaller subset (16).

867 Even though we achieved high values of predictive performances, this result was possible only thanks
 868 to a high quality and artifact free dataset. Density of observations in the dataset varies greatly across
 869 Europe and different species, leaving large areas either overrepresented or poorly covered (Fig. 2).
 870 Model uncertainty around predicted probability values depends strongly on both density and quality of
 871 observations, with areas scarcely covered or underrepresented in the datasets having higher uncertainty
 872 values or becoming source of extrapolation for the model. In our maps this can be seen in various locations
 873 in Norway (no absence data available) and Iceland (no presence or absence data available) across multiple
 874 species distribution maps, in both potential and realized distributions (Fig. 9).

875 Because our interest was to map distribution of forest species on a long time scale, we had to balance
 876 between using the best Earth Observation images and time-coverage (2000–2020). Eventually we have
 877 decided not to use any of the Sentinel-2 or similar products as they only span 2016–2020. In reality,
 878 Sentinel-1,2 images are more detailed (from 20 to 10 m resolution) and could probably help increase
 879 accuracy of the more recent years.

880 We have tried implementing predictions models that are holistic i.e. include inter-specific interactions

881 between target forest species and other species through species distribution maps; in practice this was not
882 easy to implement because consistent, long-term distributions of species i.e. distribution maps from the
883 European Atlas of Forest Tree Species (San-Miguel-Ayanz, J., de Rigo, D., Caudullo, G., Houston Durrant,
884 T., Mauri, 2016) are only available at relatively coarse spatial resolutions and not for all of the forest tree
885 species. In addition, we have not considered other valuable information that can be classified among the
886 biotic factors i.e. relative abundance/richness (Kass et al., 2020) or abundance of large individuals (Meier
887 et al., 2010) due to data unavailability for both spatial and temporal scale of the study. The prediction
888 errors are provided as standard deviation of the base learners. This approach is computationally acceptable,
889 but could be further fine-tuned to increase accuracy of errors as in van den Hoogen et al. (2020).

890 4.6 Future work and directions

891 The predictions of forest tree species we have produced are possibly the most detailed and certainly the
892 largest open datasets available to date. Although we had to innovate both Species Distribution Modelling
893 methods and our computing and coding skills, we still recognize many future areas of improvements.
894 These include:

- 895 • Fusing ALL Earth Observation data available e.g. Harmonized Landsat Sentinel-2 (HLS) and
896 eventually Sentinel 1 data to help improve predictions.
- 897 • Adding and testing hyper-spectral images (e.g. from the future hyper-spectral missions such
898 as ENMAP, see <https://www.enmap.org/>) for helping increase accuracy of predictions for
899 especially species that grow under dominant species.
- 900 • Incorporating more species data, in particular those from NFI plots would further improve the
901 species predictions at European and national scales.
- 902 • Adding more sophisticated and different Machine Learning algorithms such as Deep learning
903 techniques (Lakshminarayanan et al., 2016) and similar.
- 904 • Modeling dynamics of forest tree species not as fields, but as objects through cellular automata or
905 similar algorithms.
- 906 • Using spatio-temporal Machine Learning to predict future states of forest tree distributions using
907 climate scenarios and similar.

908 We have released the maps and the code under open data / open source licenses to enable other similar
909 research and to help speed up land restoration and reforestation projects in Europe. The code is avail-
910 able in our GitLab repository at [https://gitlab.com/geoharmonizer_inea/spatial-layers/
911 -/tree/master/veg_mapping](https://gitlab.com/geoharmonizer_inea/spatial-layers/-/tree/master/veg_mapping). To suggest any improvement visit instead [https://gitlab.com/
912 geoharmonizer_inea/spatial-layers/-/issues](https://gitlab.com/geoharmonizer_inea/spatial-layers/-/issues). The predictions of tree species are available as
913 Cloud-Optimized GeoTIFFs via <https://maps.opendatascience.eu> and can be displayed using
914 compare tool in 2D and 3D. (Fig. 11).

915 European forest dynamics, even though some recent results indicate increased mortality in forest tree
916 species e.g. Popkin (2021), is probably among the least troubling in comparison to other continents. Our



Figure 11. Difference between potential and realized distribution for *Fagus sylvatica* in Northern Spain for the period 2018–2020 visualized using slider in the Open Data Science Europe viewer (<https://maps.opendatascience.eu>).

917 methodological framework could thus potentially be implemented also at a global scale, and possibly
 918 through Google Earth Engine (van den Hoogen et al., 2021) or through the ESA’s OpenEO platform
 919 (<https://openeo.cloud/>) to produce high resolution (10–30 m) predictions of global forest dynamics.
 920 Our estimate is that globally there are much more forest tree species important for forest management
 921 and monitoring i.e. at the scale of 500–1000 m. For example, in Brazil, it has been estimated that
 922 about 220 forest tree species cover most of the land and represent over 95 of the biomass i.e. so called
 923 “hyper-dominant species” (Draper et al., 2021). Scaling up approach described in this paper to help produce
 924 objective predictions and help monitor forest dynamics and support re-forestation efforts across globe is
 925 our next frontier.

926 ACKNOWLEDGEMENTS

927 This work is co-financed by the European Union CEF Telecom project 2018-EU-IA-0095. We are
 928 grateful to the GiLAB company from Belgrade, Serbia for their support with processing and publishing
 929 produced data via the opendatascience.eu data portal. We are also grateful to the Geo-harmonizer project
 930 partners CVUT Prague, mundialis, Terrasigna & MultiOne for helping with processing all LUCAS ground
 931 observations and quality control.

932 CONFLICT OF INTEREST

933 The authors confirm no conflict of interest.

934 REFERENCES

- 935 Araújo, M. B. and New, M. (2007). Ensemble forecasting of species distributions. *Trends in Ecology &*
936 *Evolution*, 22(1):42–47.
- 937 Astola, H., Häme, T., Sirro, L., Molinier, M., and Kilpi, J. (2019). Comparison of Sentinel-2 and Landsat
938 8 imagery for forest variable prediction in boreal region. *Remote Sensing of Environment*, 223:257–273.
- 939 Bechtel, B. and Schmidt, K. J. (2011). Floristic mapping data as a proxy for the mean urban heat island.
940 *Climate Research*, 49(1):45–58.
- 941 Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., and Jones, Z. M.
942 (2016). mlr: Machine learning in r. *Journal of Machine Learning Research*, 17(170):1–5.
- 943 Booth, T. H. (2018). Species distribution modelling tools and databases to assist managing forests under
944 climate change. *Forest Ecology and Management*, 430(15):196–203.
- 945 Booth, T. H., Nix, H. A., Busby, J. R., and Hutchinson, M. F. (2014). BIOCLIM : the first species
946 distribution modelling package, its early applications and relevance to most current MaxEnt studies.
947 *Diversity and Distributions*, 20(1):1–9.
- 948 Bradley, B. A., Olsson, A. D., Wang, O., Dickson, B. G., Pelech, L., Sesnie, S. E., and Zachmann, L. J.
949 (2012). Species detection vs. habitat suitability: Are we biasing habitat suitability models with remotely
950 sensed data? *Ecological Modelling*, 244:57–64.
- 951 Breiman, L. (2001a). Random forests. *Machine learning*, 45(1):5–32.
- 952 Breiman, L. (2001b). Statistical Modeling: The Two Cultures (with comments and a rejoinder by the
953 author). *Statistical Science*, 16(3):199–231.
- 954 Brown, K. E., Bhuiyan, F. A., and Talbert, D. A. (2020). Uncertainty quantification in multimodal
955 ensembles of deep learners. In *The Thirty-Third International Flairs Conference*.
- 956 Brus, D., Hengeveld, G., Walvoort, D., Goedhart, P., Heidema, A., Nabuurs, G., and Gunia, K. (2012).
957 Statistical mapping of tree species over europe. *European Journal of Forest Research*, 131(1):145–157.
- 958 Bucklin, D. N., Basille, M., Benschoter, A. M., Brandt, L. A., Mazzotti, F. J., Romanach, S. S., Speroterra,
959 C., and Watling, J. I. (2015). Comparing species distribution models constructed with different subsets
960 of environmental predictors. *Diversity and distributions*, 21(1):23–35.
- 961 Büttner, G., Steenmans, C., Bossard, M., Feranec, J., and Kolár, J. (1998). The European CORINE land
962 cover database. *International Archives of Photogrammetry and Remote Sensing*, 32:633–638.
- 963 Cord, A., Schmidt, M., and Dech, S. (2009). Potential and limitations of multi-temporal Earth observation
964 data to improve model results of tree species distribution in Mexico. *Proceedings, 33rd International*
965 *Symposium on Remote Sensing of Environment, ISRSE 2009*, pages 285–288.
- 966 d’Andrimont, R., Verhegghen, A., Meroni, M., Lemoine, G., Strobl, P., Eiselt, B., Yordanov, M., Martinez-
967 Sanchez, L., and van der Velde, M. (2021). Lucas copernicus 2018: Earth-observation-relevant in situ
968 data on land cover and use throughout the european union. *Earth System Science Data*, 13(3):1119–
969 1133.
- 970 de Rigo, D., Caudullo, G., Houston Durrant, T., and San-Miguel-Ayanz, J. (2016a). The European Atlas
971 of Forest Tree Species: modelling, data and information on forest tree species. *European Atlas of*
972 *Forest Tree Species*, page e01aa69.

- 973 de Rigo, D., Houston Durrant, T., Caudullo, G., and Barredo, J. I. (2016b). European forests: an ecological
974 overview. In *European Atlas of Forest Tree Species*, pages 24 – 31. Publication Office of the European
975 Union, Luxembourg.
- 976 Domke, G. M., Oswalt, S. N., Walters, B. F., and Morin, R. S. (2020). Tree planting has the potential
977 to increase carbon sequestration capacity of forests in the united states. *Proceedings of the National
978 Academy of Sciences*, 117(40):24649–24651.
- 979 Dottori, F., Salamon, P., Bianchi, A., Alfieri, L., Hirpa, F. A., and Feyen, L. (2016). Development and
980 evaluation of a framework for global flood hazard mapping. *Advances in Water Resources*, 94:87–102.
- 981 Draper, F. C., Costa, F. R., Arellano, G., Phillips, O. L., Duque, A., Macía, M. J., Ter Steege, H., Asner,
982 G. P., Berenguer, E., Schiatti, J., et al. (2021). Amazon tree dominance across forest strata. *Nature
983 ecology & evolution*, 5(6):757–767.
- 984 Easdale, M. H., Bruzzone, O., Mapfumo, P., and Tittonell, P. (2018). Phases or regimes? R evisiting
985 NDVI trends as proxies for land degradation. *Land Degradation & Development*, 29(3):433–445.
- 986 Elith, J., Leathwick, J. R., and Hastie, T. (2008). A working guide to boosted regression trees. *Journal of
987 Animal Ecology*, 77(4):802–813.
- 988 EUROSTAT (2017). Land Cover/Use Statistics (LUCAS) Database.
- 989 Fawcett, D., Bennie, J., and Anderson, K. (2021). Monitoring spring phenology of individual tree crowns
990 using drone-acquired ndvi data. *Remote Sensing in Ecology and Conservation*, 7(2):227–244.
- 991 Ferri, C., Hernández-Orallo, J., and Modroui, R. (2009). An experimental comparison of performance
992 measures for classification. *Pattern Recognition Letters*, 30(1):27–38.
- 993 Fick, S. E. and Hijmans, R. J. (2017). WorldClim 2: new 1-km spatial resolution climate surfaces for
994 global land areas. *International journal of climatology*, 37(12):4302–4315.
- 995 Fix, E. and Hodges, J. L. (1989). Discriminatory Analysis. Nonparametric Discrimination: Consistency
996 Properties. *International Statistical Review / Revue Internationale de Statistique*, 57(3):238.
- 997 Fois, M., Cuenca-Lombraña, A., Fenu, G., and Bacchetta, G. (2018). Using species distribution models
998 at local scale to guide the search of poorly known species: Review, methodological issues and future
999 directions. *Ecological Modelling*, 385:124–132.
- 1000 Fourcade, Y., Besnard, A. G., and Secondi, J. (2018). Paintings predict the distribution of species, or
1001 the challenge of selecting environmental predictors and evaluation statistics. *Global Ecology and
1002 Biogeography*, 27(2):245–256.
- 1003 Franklin, J. (2010). *Mapping Species Distributions: Spatial Inference and Prediction*. Cambridge
1004 University Press, Cambridge.
- 1005 Friedman, J. H. (2002). Stochastic gradient boosting. *Computational Statistics & Data Analysis*,
1006 38(4):367–378. Nonlinear Methods and Data Mining.
- 1007 Gao, B.-C. (1996). NDWI—A normalized difference water index for remote sensing of vegetation liquid
1008 water from space. *Remote sensing of environment*, 58(3):257–266.
- 1009 Gao, T., Zhu, J., Zheng, X., Shang, G., Huang, L., and Wu, S. (2015). Mapping spatial distribution of
1010 larch plantations from multi-seasonal landsat-8 oli imagery and multi-scale textures using random
1011 forests. *Remote Sensing*, 7(2):1702–1720.
- 1012 Gelfand, A. E. and Shirota, S. (2021). The role of odds ratios in joint species distribution modeling.

- 1013 *Environmental and Ecological Statistics*, 28(2):287–302.
- 1014 Gobeyn, S., Mouton, A. M., Cord, A. F., Kaim, A., Volk, M., and Goethals, P. L. (2019). Evolutionary
1015 algorithms for species distribution modelling: A review in the context of machine learning. *Ecological*
1016 *Modelling*, 392(June 2018):179–195.
- 1017 Gomes, C., Nocairi, H., Thomas, M., Ibanez, F., Collin, J.-F., and Saporta, G. (2012). Stacking prediction
1018 for a binary outcome. In *Compstat 2012*, pages 271–282, Limassol, Cyprus.
- 1019 Gorelick, N., Hancher, M., Dixon, M., Ilyushchenko, S., Thau, D., and Moore, R. (2017). Google Earth
1020 Engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment*, 202:18–27.
- 1021 Gottschalk, T. K., Aue, B., Hotes, S., and Ekschmitt, K. (2011). Influence of grain size on species-habitat
1022 models. *Ecological Modelling*, 222(18):3403–3412.
- 1023 Guisan, A., Petitpierre, B., Broennimann, O., Daehler, C., and Kueffer, C. (2014). Unifying niche shift
1024 studies: Insights from biological invasions. *Trends in Ecology and Evolution*, 29(5):260–269.
- 1025 Guisan, A. and Thuiller, W. (2005). Predicting species distribution: Offering more than simple habitat
1026 models. *Ecology Letters*, 8(9):993–1009.
- 1027 Guisan, A., Thuiller, W., and Zimmermann, N. E. (2017). *Habitat Suitability and Distribution Models*.
1028 Cambridge University Press, Cambridge.
- 1029 Guisan, A., Tingley, R., Baumgartner, J. B., Naujokaitis-Lewis, I., Sutcliffe, P. R., Tulloch, A. I., Regan,
1030 T. J., Brotons, L., McDonald-Madden, E., Mantyka-Pringle, C., Martin, T. G., Rhodes, J. R., Maggini,
1031 R., Setterfield, S. A., Elith, J., Schwartz, M. W., Wintle, B. A., Broennimann, O., Austin, M., Ferrier,
1032 S., Kearney, M. R., Possingham, H. P., and Buckley, Y. M. (2013). Predicting species distributions for
1033 conservation decisions. *Ecology Letters*, 16(12):1424–1435.
- 1034 Hampton, S. E., Strasser, C. A., Tewksbury, J. J., Gram, W. K., Budden, A. E., Batcheller, A. L., Duke,
1035 C. S., and Porter, J. H. (2013). Big data and the future of ecology. *Frontiers in Ecology and the*
1036 *Environment*, 11(3):156–162.
- 1037 Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A., Thau, D.,
1038 Stehman, S. V., Goetz, S. J., Loveland, T. R., Kommareddy, A., Egorov, A., Chini, L., Justice, C. O.,
1039 and Townshend, J. R. G. (2013). High-resolution global maps of 21st-century forest cover change.
1040 *Science*, 342(6160):850–853.
- 1041 Hao, T., Elith, J., Guillera-Aroita, G., and Lahoz-Monfort, J. J. (2019). A review of evidence about
1042 use and performance of species distribution modelling ensembles like BIOMOD. *Diversity and*
1043 *Distributions*, 25(5):839–852.
- 1044 Hao, T., Elith, J., Lahoz-Monfort, J. J., and Guillera-Aroita, G. (2020). Testing whether ensemble
1045 modelling is advantageous for maximising predictive performance of species distribution models.
1046 *Ecography*, 43(4):549–558.
- 1047 Hastie, T., Qian, J., and Tay, K. (2016). An introduction to glmnet.
- 1048 He, K. S., Zhang, J., and Zhang, Q. (2009). Linking variability in species composition and MODIS NDVI
1049 based on beta diversity measurements. *acta oecologica*, 35(1):14–21.
- 1050 Hefley, T. J. and Hooten, M. B. (2016). Hierarchical species distribution models. *Current Landscape*
1051 *Ecology Reports*, 1(2):87–97.
- 1052 Heisig, J. and Hengl, T. (2020). Harmonized Tree Species Occurrence Points for Europe. URL:

- 1053 <https://zenodo.org/record/4061816>, Dataset Version: 0.2.
- 1054 Hengl, T., Leal Parente, L., Krizan, J., and Bonannella, C. (2020). Continental Europe Digital Terrain
1055 Model at 30 m resolution based on GEDI, ICESat-2, AW3D, GLO-30, EUDEM, MERIT DEM and
1056 background layers.
- 1057 Hengl, T., Walsh, M. G., Sanderman, J., Wheeler, I., Harrison, S. P., and Prentice, I. C. (2018). Global
1058 mapping of potential natural vegetation: an assessment of machine learning algorithms for estimating
1059 land potential. *PeerJ*, 6:e5457.
- 1060 Hierro, J. L., Maron, J. L., and Callaway, R. M. (2005). A biogeographical approach to plant invasions:
1061 the importance of studying exotics in their introduced and native range. *Journal of ecology*, 93(1):5–15.
- 1062 Hijmans, R. J., Cameron, S. E., Parra, J. L., Jones, P. G., and Jarvis, A. (2005). Very high resolution
1063 interpolated climate surfaces for global land areas. *International Journal of Climatology: A Journal of
1064 the Royal Meteorological Society*, 25(15):1965–1978.
- 1065 Hill, L., Hector, A., Hemery, G., Smart, S., Tanadini, M., and Brown, N. (2017). Abundance distributions
1066 for tree species in Great Britain: A two-stage approach to modeling abundance using species distribution
1067 modeling and random forest. *Ecology and evolution*, 7(4):1043–1056.
- 1068 Hoffer, R. (1984). Remote sensing to measure the distribution and structure of vegetation. *The Role of
1069 Terrestrial Vegetation in the Global Carbon Cycle: Measurement by Remote Sensing*, pages 131–59.
- 1070 Huete, A., Didan, K., Miura, T., Rodriguez, E. P., Gao, X., and Ferreira, L. G. (2002). Overview of
1071 the radiometric and biophysical performance of the modis vegetation indices. *Remote sensing of
1072 environment*, 83(1-2):195–213.
- 1073 Huete, A. R. (1988). A soil-adjusted vegetation index (SAVI). *Remote sensing of environment*, 25(3):295–
1074 309.
- 1075 Hutchinson, G. E. (1957). Concluding Remarks. *Cold Spring Harbor Symposia on Quantitative Biology*,
1076 22:415–427.
- 1077 Iturbide, M., Bedia, J., and Gutiérrez, J. M. (2018). Background sampling and transferability of species
1078 distribution model ensembles under climate change. *Global and Planetary Change*, 166(March):19–29.
- 1079 Iturbide, M., Bedia, J., Herrera, S., del Hierro, O., Pinto, M., and Gutiérrez, J. M. (2015). A framework
1080 for species distribution modelling with improved pseudo-absence generation. *Ecological Modelling*,
1081 312:166–174.
- 1082 Jiang, Z., Huete, A. R., Didan, K., and Miura, T. (2008). Development of a two-band enhanced vegetation
1083 index without a blue band. *Remote sensing of Environment*, 112(10):3833–3845.
- 1084 Jiménez-Valverde, A., Lobo, J. M., and Hortal, J. (2008). Not as good as they seem: The importance of
1085 concepts in species distribution modelling. *Diversity and Distributions*, 14(6):885–890.
- 1086 Jiménez-Valverde, A., Peterson, A. T., Soberón, J., Overton, J., Aragón, P., and Lobo, J. M. (2011). Use
1087 of niche models in invasive species risk assessments. *Biological invasions*, 13(12):2785–2797.
- 1088 Kaky, E., Nolan, V., Alatawi, A., and Gilbert, F. (2020). A comparison between Ensemble and MaxEnt
1089 species distribution modelling approaches for conservation: A case study with Egyptian medicinal
1090 plants. *Ecological Informatics*, 60:101–150.
- 1091 Karger, D. N., Conrad, O., Böhner, J., Kawohl, T., Kreft, H., Soria-Auza, R. W., Zimmermann, N. E.,
1092 Linder, H. P., and Kessler, M. (2017). Climatologies at high resolution for the earth's land surface areas.

- 1093 *Scientific data*, 4(1):1–20.
- 1094 Karger, D. N., Dabaghchian, B., Lange, S., Thuiller, W., Zimmermann, N. E., and Graham, C. H. (2020).
1095 High resolution climate data for europe.
- 1096 Kass, J. M., Anderson, R. P., Espinosa-Lucas, A., Juárez-Jaimes, V., Martínez-Salas, E., Botello, F.,
1097 Tavera, G., Flores-Martínez, J. J., and Sánchez-Cordero, V. (2020). Biotic predictors with phenolog-
1098 ical information improve range estimates for migrating monarch butterflies in mexico. *Ecography*,
1099 43(3):341–352.
- 1100 Keenan, R. J. (2015). Climate change impacts and adaptation in forest management: a review. *Annals of*
1101 *Forest Science*, 72(2):145–167.
- 1102 Key, C. H. and Benson, N. C. (1999). The Normalized Burn Ratio (NBR): A Landsat TM radiometric
1103 measure of burn severity. *United States Geological Survey, Northern Rocky Mountain Science Center:*
1104 *Bozeman, MT, USA*.
- 1105 Kindt, R. (2018). Ensemble species distribution modelling with transformed suitability values. *Environ-*
1106 *mental Modelling & Software*, 100:136–145.
- 1107 Lakshminarayanan, B., Pritzel, A., and Blundell, C. (2016). Simple and scalable predictive uncertainty
1108 estimation using deep ensembles. *arXiv preprint arXiv:1612.01474*.
- 1109 Lefebvre, D., Williams, A. G., Kirk, G. J., Burgess, J., Meersmans, J., Silman, M. R., Román-Dañobeytia,
1110 F., Farfan, J., Smith, P., et al. (2021). Assessing the carbon capture potential of a reforestation project.
1111 *Scientific reports*, 11(1):1–10.
- 1112 Madonsela, S., Cho, M. A., Ramoelo, A., and Mutanga, O. (2017). Remote sensing of species diversity
1113 using Landsat 8 spectral variables. *ISPRS Journal of Photogrammetry and Remote Sensing*, 133:116–
1114 127.
- 1115 Manzoor, S. A., Griffiths, G., and Lukac, M. (2018). Species distribution model transferability and model
1116 grain size-finer may not always be better. *Scientific Reports*, 8(1):1–9.
- 1117 Marchi, M. and Ducci, F. (2018). Some refinements on species distribution models using tree-level
1118 National Forest Inventories for supporting forest management and marginal forest population detection.
1119 *iForest - Biogeosciences and Forestry*, 11(2):291–299.
- 1120 Martinez-Minaya, J., Cameletti, M., Conesa, D., and Pennino, M. G. (2018). Species distribution
1121 modeling: a statistical review with focus in spatio-temporal issues. *Stochastic environmental research*
1122 *and risk assessment*, 32(11):3227–3244.
- 1123 Mather, J. R. and Yoshioka, G. A. (1968). The role of climate in the distribution of vegetation. *Annals of*
1124 *the Association of American Geographers*, 58(1):29–41.
- 1125 Mauri, A., Strona, G., and San-Miguel-Ayanz, J. (2017). EU-Forest, a high-resolution tree occurrence
1126 dataset for Europe. *Scientific Data*, 4(1):160123.
- 1127 McCulloch, W. S. and Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The*
1128 *bulletin of mathematical biophysics*, 5(4):115–133.
- 1129 Mehra, A., Tripathy, P., Faridi, A., and Chinmay, A. (2019). Ensemble learning approach to improve
1130 existing models. *International Journal of Innovative Science and Research Technology*, 4.
- 1131 Meier, E. S., Kienast, F., Pearman, P. B., Svenning, J. C., Thuiller, W., Araújo, M. B., Guisan, A., and
1132 Zimmermann, N. E. (2010). Biotic and abiotic variables show little redundancy in explaining tree

- 1133 species distributions. *Ecography*, 33(6):1038–1048.
- 1134 Miller, J., Franklin, J., and Aspinall, R. (2007). Incorporating spatial dependence in predictive vegetation
1135 models. *Ecological Modelling*, 202(3-4):225–242.
- 1136 Molnar, C. (2020). *Interpretable Machine Learning*. Lulu.com.
- 1137 Nandy, S., Singh, R., Ghosh, S., Watham, T., Kushwaha, S. P. S., Kumar, A. S., and Dadhwal, V. K. (2017).
1138 Neural network-based modelling for forest biomass assessment. *Carbon Management*, 8(4):305–317.
- 1139 Nave, L. E., Walters, B. F., Hofmeister, K., Perry, C. H., Mishra, U., Domke, G. M., and Swanston, C.
1140 (2019). The role of reforestation in carbon sequestration. *New Forests*, 50(1):115–137.
- 1141 Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized linear models. *Journal of the Royal*
1142 *Statistical Society. Series A (General)*, 135(3):370–384.
- 1143 Olaya, V. (2009). Chapter 6 basic land-surface parameters. In Hengl, T. and Reuter, H. I., editors,
1144 *Geomorphometry*, volume 33 of *Developments in Soil Science*, pages 141–169. Elsevier.
- 1145 Palankar, M. R., Iamnitich, A., Ripeanu, M., and Garfinkel, S. (2008). Amazon S3 for science grids:
1146 a viable solution? In *Proceedings of the 2008 international workshop on Data-aware distributed*
1147 *computing*, pages 55–64.
- 1148 Pearson, R. G. and Dawson, T. P. (2003). Predicting the impacts of climate change on the distribution of
1149 species: are bioclimate envelope models useful? *Global ecology and biogeography*, 12(5):361–371.
- 1150 Pearson, R. G., Thuiller, W., Araújo, M. B., Martinez-Meyer, E., Brotons, L., McClean, C., Miles,
1151 L., Segurado, P., Dawson, T. P., and Lees, D. C. (2006). Model-based uncertainty in species range
1152 prediction. *Journal of biogeography*, 33(10):1704–1711.
- 1153 Pecchi, M., Marchi, M., Giannetti, F., Moriondo, M., Bernetti, I., Bindi, M., and Chirici, G. (2019).
1154 Species distribution modelling to support forest management. A literature review. *Ecological Modelling*,
1155 411(submitted):108817.
- 1156 Pekel, J.-F., Cottam, A., Gorelick, N., and Belward, A. S. (2016). High-resolution mapping of global
1157 surface water and its long-term changes. *Nature*, 540(7633):418–422.
- 1158 Pérez Chaves, P., Ruokolainen, K., and Tuomisto, H. (2018). Using remote sensing to model tree species
1159 distribution in Peruvian lowland Amazonia. *Biotropica*, 50(5):758–767.
- 1160 Peterson, A. T., Soberón, J. M., Pearson, R. G., Anderson, R. P., Martínez-Meyer, E., Nakamura, M., and
1161 Araújo, M. B. (2011). *Ecological niches and geographic distributions (MPB-49)*. Princeton University
1162 Press.
- 1163 Phillips, S. J., Dudik, M., Elith, J., Graham, C. H., Leathwick, J., Ferrier, S., Applications, S. E., Jan, N.,
1164 Phillips, J., Leathwick, J., and Elith, J. (2009). Sample Selection Bias and Presence-Only Distribution
1165 Models : Implications for Background and Pseudo-Absence Data Published by : Ecological Society of
1166 America Linked references are available on JSTOR for this article : Your use of the JSTOR archive
1167 indica. *Ecological Applications*, 19(1):181–197.
- 1168 Pigeon, G., Legain, D., Durand, P., and Masson, V. (2007). Anthropogenic heat release in an old European
1169 agglomeration (Toulouse, France). *International Journal of Climatology: A Journal of the Royal*
1170 *Meteorological Society*, 27(14):1969–1981.
- 1171 Popkin, G. (2021). Germany’s trees are dying. a fierce debate has broken out over how to respond. *Science*,
1172 374.

- 1173 Porfirio, L. L., Harris, R. M., Lefroy, E. C., Hugh, S., Gould, S. F., Lee, G., Bindoff, N. L., and Mackey,
1174 B. (2014). Improving the use of species distribution models in conservation planning and management
1175 under climate change. *PLoS ONE*, 9(11):1–21.
- 1176 Potapov, P., Hansen, M. C., Kommareddy, I., Kommareddy, A., Turubanova, S., Pickens, A., Adusei, B.,
1177 Tyukavina, A., and Ying, Q. (2020). Landsat analysis ready data for global land cover and land cover
1178 change mapping. *Remote Sensing*, 12(3):426.
- 1179 Prasad, A. M., Iverson, L. R., and Liaw, A. (2006). Newer classification and regression tree techniques:
1180 Bagging and random forests for ecological prediction. *Ecosystems*, 9(2):181–199.
- 1181 Prates-Clark, C. D. C., Saatchi, S. S., and Agosti, D. (2008). Predicting geographical distribution models
1182 of high-value timber trees in the Amazon Basin using remotely sensed data. *Ecological Modelling*,
1183 211(3-4):309–323.
- 1184 Qi, J., Chehbouni, A., Huete, A. R., Kerr, Y. H., and Sorooshian, S. (1994). A modified soil adjusted
1185 vegetation index. *Remote sensing of environment*, 48(2):119–126.
- 1186 Quinlan, J. R. (1986). Induction of decision trees. *Machine Learning*, 1(1):81–106.
- 1187 Ripley, B. and Venables, W. (2017). *nnet: Feed-Forward Neural Networks and Multinomial Log-Linear*
1188 *Models*. R package version 7.3-12.
- 1189 Roberts, D. R., Bahn, V., Ciuti, S., Boyce, M. S., Elith, J., Guillera-Arroita, G., Hauenstein, S., Lahoz-
1190 Monfort, J. J., Schröder, B., Thuiller, W., et al. (2017). Cross-validation strategies for data with
1191 temporal, spatial, hierarchical, or phylogenetic structure. *Ecography*, 40(8):913–929.
- 1192 San-Miguel-Ayanz, J., de Rigo, D., Caudullo, G., Houston Durrant, T., Mauri, A. E. (2016). *European*
1193 *Atlas of Forest Tree Species*. Publication Office of the European Union.
- 1194 Schloss, A., Kicklighter, D., Kaduk, J., Wittenberg, U., and Intercomparison, T. P. O. T. P. N. M. (1999).
1195 Comparing global models of terrestrial net primary productivity (NPP): comparison of NPP to climate
1196 and the Normalized Difference Vegetation Index (NDVI). *Global Change Biology*, 5(S1):25–34.
- 1197 Schratz, P., Muenchow, J., Iturriza, E., Richter, J., and Brenning, A. (2019). Hyperparameter tuning and
1198 performance assessment of statistical and machine-learning algorithms using spatial data. *Ecological*
1199 *Modelling*, 406:109–120.
- 1200 Shabani, F., Kumar, L., and Ahmadi, M. (2018). Assessing accuracy methods of species distribution
1201 models: AUC, Specificity, Sensitivity and the True Skill Statistic. *Global Journal of Human Social*
1202 *Science*, 18(1):6–18.
- 1203 Shi, X., Wong, Y. D., Li, M. Z.-F., Palanisamy, C., and Chai, C. (2019). A feature learning approach based
1204 on XGBoost for driving assessment and risk prediction. *Accident Analysis & Prevention*, 129:170–179.
- 1205 Soberón, J. M. (2010). Niche and area of distribution modeling: A population ecology perspective.
1206 *Ecography*, 33(1):159–167.
- 1207 Soberón, J. M. and Peterson, A. T. (2005). Interpretation of Models of Fundamental Ecological Niches
1208 and Species' Distributional Areas. *Biodiversity Informatics*, 2(10):1–10.
- 1209 Strickland, G. E. I., Luther, J. E., White, J. C., Wulder, M. A., Strickland, G. E. I., Luther, J. E., White, J. C.,
1210 and Wulder, M. A. (2020). Extending Estimates of Tree and Tree Species Presence-Absence through
1211 Space and Time Using Landsat Composites. *Canadian Journal of Remote Sensing*, 46(5):567–584.
- 1212 Therneau, T. M. and Atkinson, E. J. (2011). An Introduction to Recursive Partitioning Using the RPART

1213 Routines. *Mayo clinic*, 61:33.

1214 Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*
1215 *Society: Series B (Methodological)*, 58(1):267–288.

1216 Tucker, C. J. (1979). Red and photographic infrared linear combinations for monitoring vegetation.
1217 *Remote sensing of Environment*, 8(2):127–150.

1218 Valavi, R., Guillera-Arroita, G., Lahoz-Monfort, J. J., and Elith, J. (2021). Predictive performance of
1219 presence-only species distribution models: a benchmark study with reproducible code. *Ecological*
1220 *Monographs*, 0(0):1–27.

1221 van den Hoogen, J., Geisen, S., Wall, D. H., Wardle, D. A., Traunspurger, W., de Goede, R. G., Adams,
1222 B. J., Ahmad, W., Ferris, H., Bardgett, R. D., et al. (2020). A global database of soil nematode
1223 abundance and functional group composition. *Scientific data*, 7(1):1–8.

1224 van den Hoogen, J., Robmann, N., Routh, D., Lauber, T., van Tiel, N., Danylo, O., and Crowther, T. W.
1225 (2021). A geospatial mapping pipeline for ecologists. *BioRxiv*.

1226 Walthert, L. and Meier, E. S. (2017). Tree species distribution in temperate forests is more influenced by
1227 soil than by climate. *Ecology and Evolution*, 7(22):9473–9484.

1228 Wang, H., Liu, H., Huang, N., Bi, J., Ma, X., Ma, Z., Shanguan, Z., Zhao, H., Feng, Q., Liang, T., et al.
1229 (2021). Satellite-derived NDVI underestimates the advancement of alpine vegetation growth over the
1230 past three decades.

1231 Weigel, R., Gilles, J., Klisz, M., Manthey, M., and Kreyling, J. (2019). Forest understory vegetation is
1232 more related to soil than to climate towards the cold distribution margin of european beech. *Journal of*
1233 *Vegetation Science*, 30(4):746–755.

1234 Wilson, A. M. and Jetz, W. (2016). Remotely sensed high-resolution global cloud dynamics for predicting
1235 ecosystem and biodiversity distributions. *PLoS biology*, 14(3):e1002415.

1236 Witjes, M., Parente, L., van Diemen, C. J., Hengl, T., Landa, M., Brodsky, L., Halounova, L., Krizan,
1237 J., Antonic, L., Ilie, C. M., et al. (2021). A spatiotemporal ensemble machine learning framework for
1238 generating land use/land cover time-series maps for Europe (2000–2019) based on LUCAS, CORINE
1239 and GLAD Landsat. *PeerJ*.

1240 Wolpert, D. H. (1992). Stacked generalization. *Neural networks*, 5(2):241–259.

1241 Zhang, C. and Ma, Y. (2012). *Ensemble machine learning: methods and applications*. Springer.

1242 Zhou, Z.-H. (2019). *Ensemble methods: foundations and algorithms*. Chapman and Hall/CRC.

1243 Zhu, Z., Wulder, M. A., Roy, D. P., Woodcock, C. E., Hansen, M. C., Radeloff, V. C., Healey, S. P., Schaaf,
1244 C., Hostert, P., Strobl, P., et al. (2019). Benefits of the free and open Landsat data policy. *Remote*
1245 *Sensing of Environment*, 224:382–385.