

A peer-reviewed version of this preprint was published in PeerJ on 22 August 2018.

[View the peer-reviewed version](https://peerj.com/articles/5457) (peerj.com/articles/5457), which is the preferred citable publication unless you specifically need to cite this preprint.

Hengl T, Walsh MG, Sanderman J, Wheeler I, Harrison SP, Prentice IC. 2018. Global mapping of potential natural vegetation: an assessment of machine learning algorithms for estimating land potential. PeerJ 6:e5457 <https://doi.org/10.7717/peerj.5457>

Global mapping of potential natural vegetation: an assessment of Machine Learning algorithms for estimating land potential

Tomislav Hengl^{Corresp., 1}, Markus G Walsh^{2,3}, Jonathan Sanderman⁴, Ichsani Wheeler¹, Sandy P Harrison⁵, Iain C Prentice⁶

¹ Envirometrix Ltd, Wageningen, Netherlands

² The Earth Institute, Columbia University, New York, United States

³ Selian Agricultural Research Inst., Arusha, Tanzania

⁴ Woods Hole Research Center, Falmouth, United States

⁵ School of Archeology, Geography and Environmental Science, University of Reading, Reading, United Kingdom

⁶ Department of Life Sciences and Grantham Institute - Climate Change and the Environment, Imperial College London, London, United Kingdom

Corresponding Author: Tomislav Hengl

Email address: tom.hengl@envirometrix.net

Potential Natural Vegetation (PNV) is the vegetation cover in equilibrium with climate, that would exist at a given location if not impacted by human activities. PNV is useful for raising public awareness about land degradation and for estimating land potential. This paper presents results of assessing Machine Learning Algorithms (MLA) — neural networks (nnet package), random forest (ranger), gradient boosting (gmb), K-nearest neighborhood (class) and cubist — for operational mapping of PNV. Three case studies were considered: (1) global distribution of biomes based on the BIOME 6000 data set (8057 modern pollen-based site reconstructions), (2) distribution of forest tree taxa in Europe based on detailed occurrence records (1,546,435 ground observations), and (3) global monthly Fraction of Absorbed Photosynthetically Active Radiation (FAPAR) values (30,301 randomly-sampled points). A stack of 160 global maps representing biophysical conditions over land, including atmospheric, climatic, relief and lithologic variables, were used as explanatory variables. The overall results indicate that random forest gives the overall best performance. The highest accuracy for predicting BIOME 6000 classes (20) was estimated to be between 33% (with spatial Cross Validation) and 68% (simple random subsetting), with the most important predictors being total annual precipitation, monthly temperatures and bioclimatic layers. Predicting forest tree species (73) resulted in mapping accuracy of 25%, with the most important predictors being monthly cloud fraction, mean annual and monthly temperatures and elevation. Regression models for FAPAR (monthly images) gave an R-square of 90% with the most important predictors being total annual precipitation, monthly cloud fraction, CHELSA bioclimatic layers and month of the year, respectively.

Further developments of PNV mapping could include using all GBIF records to map the global distribution of plant species at different taxonomic levels. This methodology could also be extended to dynamic modeling of PNV, so that future climate scenarios can be incorporated. Global maps of biomes, FAPAR and tree species at 1 km spatial resolution are available for download via <http://dx.doi.org/10.7910/DVN/QQHCIK>.

1 **Global Mapping of Potential Natural**
2 **Vegetation: An Assessment of Machine**
3 **Learning Algorithms for Estimating Land**
4 **Potential**

5 **Tomislav Hengl¹, Markus G. Walsh^{2,3}, Jonathan Sanderman⁴, Ichsani**
6 **Wheeler¹, Sandy P. Harrison⁵, and Iain C. Prentice⁶**

7 ¹**Envirometrix Ltd., Wageningen, the Netherlands**

8 ²**The Earth Institute, Columbia University, USA**

9 ³**Selian Agricultural Research Inst., Arusha, Tanzania**

10 ⁴**Woods Hole Research Center, MA USA**

11 ⁵**School of Archeology, Geography and Environmental Science, University of Reading,**
12 **UK**

13 ⁶**AXA Chair of Biosphere and Climate Impacts, Grand Challenges in Ecosystem and the**
14 **Environment, Department of Life Sciences and Grantham Institute — Climate Change**
15 **and the Environment, Imperial College London, UK**

16 Corresponding author:

17 Tomislav Hengl¹

18 Email address: tom.hengl@envirometrix.net

19 **ABSTRACT**

20 Potential Natural Vegetation (PNV) is the vegetation cover in equilibrium with climate, that would exist at
21 a given location if not impacted by human activities. PNV is useful for raising public awareness about land
22 degradation and for estimating land potential. This paper presents results of assessing Machine Learning
23 Algorithms (MLA) — neural networks (nnet package), random forest (ranger), gradient boosting (gmb),
24 K-nearest neighborhood (class) and cubist — for operational mapping of PNV. Three case studies were
25 considered: (1) global distribution of biomes based on the BIOME 6000 data set (8057 modern pollen-based
26 site reconstructions), (2) distribution of forest tree taxa in Europe based on detailed occurrence records
27 (1,546,435 ground observations), and (3) global monthly Fraction of Absorbed Photosynthetically Active
28 Radiation (FAPAR) values (30,301 randomly-sampled points). A stack of 160 global maps representing
29 biophysical conditions over land, including atmospheric, climatic, relief and lithologic variables, were used
30 as explanatory variables. Overall, random forest models gave the best performance. The highest accuracy
31 for predicting BIOME 6000 classes (20) was estimated to be between 33 % (with spatial Cross Validation)
32 and 68 % (simple random subsetting), with the most important predictors being total annual precipitation,
33 monthly temperatures and bioclimatic layers. Predicting forest tree species (73) resulted in mapping
34 accuracy of 25 %, with the most important predictors being monthly cloud fraction, mean annual and
35 monthly temperatures and elevation. Regression models for FAPAR (monthly images) gave an R-square of
36 90 % with the most important predictors being total annual precipitation, monthly cloud fraction, CHELSA
37 bioclimatic layers and month of the year, respectively. Further developments of PNV mapping could include
38 using all GBIF records to map the global distribution of plant species at different taxonomic levels. This
39 methodology could also be extended to dynamic modeling of PNV, so that future climate scenarios can be
40 incorporated. Global maps of biomes, FAPAR and tree species at 1 km spatial resolution are available for
41 download via <http://dx.doi.org/10.7910/DVN/QQHCIK>.

42 Submitted to PeerJ on 26th of March 2018; 1st revision on 8th of July 2018;

43 INTRODUCTION

44 Potential Natural Vegetation (PNV) is the “*vegetation cover in equilibrium with climate, that would exist*
45 *at a given location non-impacted by human activities*” (Levvasseur et al., 2012; Østbye Hemsing and
46 Bryn, 2012). It is a hypothetical vegetation state assuming natural (undisturbed) physical conditions, a
47 reference status of vegetation assuming no degradation and/or no unusual ecological disturbances. PNV is
48 especially useful for raising public awareness about land degradation (Weisman, 2012) and for estimating
49 land potential (Herrick et al., 2013). For example, Omernik (1987) details PNV maps for USA; Bohn et al.
50 (2007) provides maps for EU; Carnahan (1989) for Australia; Marinova et al. (2018) maps PNV for the
51 Eastern Mediterranean–Black Sea–Caspian–Corridor; and maps of PNV for Latin America are available
52 in Marchant et al. (2009). Regarding specific tree species, San-Miguel-Ayanz et al. (2016) provide habitat
53 suitability maps for the main forest tree species in Europe, based on environmental variables, especially
54 bioclimatic variables such as average temperature of the coldest month, precipitation of the driest month
55 and similar. Potapov et al. (2011) generated a global map of potential forest cover at 1 km resolution
56 (publicly available from <http://globalforestwatch.org/map/>). Erb et al. (2017) produced a global
57 map of potential biomass stocks by reversing the current managed land use systems to natural vegetation.
58 Levvasseur et al. (2012) and Tian et al. (2016) predict global PNV classes using environmental covariates

59 such as climatic images and landform parameters. [Griscom et al. \(2017\)](#) recently produced a global
60 reforestation map at 1 km resolution.

61 A common limitation of existing maps is their coarse spatial resolution (about 25 km) limiting the
62 use of these maps for operational planning (e.g. [Marchant et al. \(2009\)](#); [Levvasseur et al. \(2012\)](#) and
63 [Tian et al. \(2016\)](#)). In addition, comparisons of multiple overlapping sources of PNV maps shows that
64 they rarely agree with one another since they do not share the same mapping criteria and, traditionally,
65 emphasize regionally-specific botanical groupings rather than functional classifications. Limitations of
66 maps based on field surveys of PNV (e.g., [Bohn et al. \(2007\)](#)) arise from assumptions about controls on
67 vegetation distribution based on extrapolation from a limited number of field surveys.

68 Here we provide an update of comparable global PNV maps produced by [Potapov et al. \(2011\)](#);
69 [Levvasseur et al. \(2012\)](#); [Tian et al. \(2016\)](#) and [Erb et al. \(2017\)](#). We explore the possibility of increasing
70 the mapping accuracy using up-to-date maps of climate, atmosphere dynamics, landform and lithology,
71 and state-of-the-art machine learning methods. Our final aim is to produce PNV maps that are more
72 detailed, richer in information, based on objective reproducible methods; and potentially more usable
73 for global modeling and awareness raising projects. We focus on improving the spatial detail, thematic
74 accuracy and reproducibility of maps, at the cost of increasing the total computing load. We also consider
75 automation of the prediction process so that the maps can be rapidly updated as new ground truth data is
76 obtained. Our modeling follows three phases:

- 77 (a) model selection: we compare possible models of interest for PNV mapping and choose the optimal
78 spatial prediction framework based on the cross-validation results,
- 79 (b) model assessment: we assess the uncertainty of predictions per vegetation class and try to determine
80 objectively the limitations of the mapping products for wider uses, and
- 81 (c) prediction: we use the best performing models to produce spatial predictions, then visually assess
82 maps and if necessary repeat steps a–c.

83 **METHODS AND MATERIALS**

84 **Theory**

85 PNV is the hypothetical vegetation cover that would be present if the vegetation were in equilibrium with
86 environmental controls, including climatic factors and disturbance, and not subject to human management.
87 When considering PNV, one needs to distinguish between potential “*natural*” and potential “*managed*”
88 vegetation, and “*actual*” natural and “*actual*” managed vegetation (Fig. 1a). Vegetation is in general
89 a dynamic feature. Also PNV changes as the climatic conditions change. For example, with the future
90 global warming and changes in our climate, PNV might be significantly different than pre-industrial
91 revolution. Therefore it is important to reference PNV to the time period of interest, so that historic PNV
92 and current or future PNV maps can be produced (Fig. 1b).

93 In addition to the differentiation between the potential and actual natural vegetation, there are also
94 three sub-types of the PNV that need to be considered:

- 95 1. PNV model A: based on the autochthonous or native vegetation and living species only.

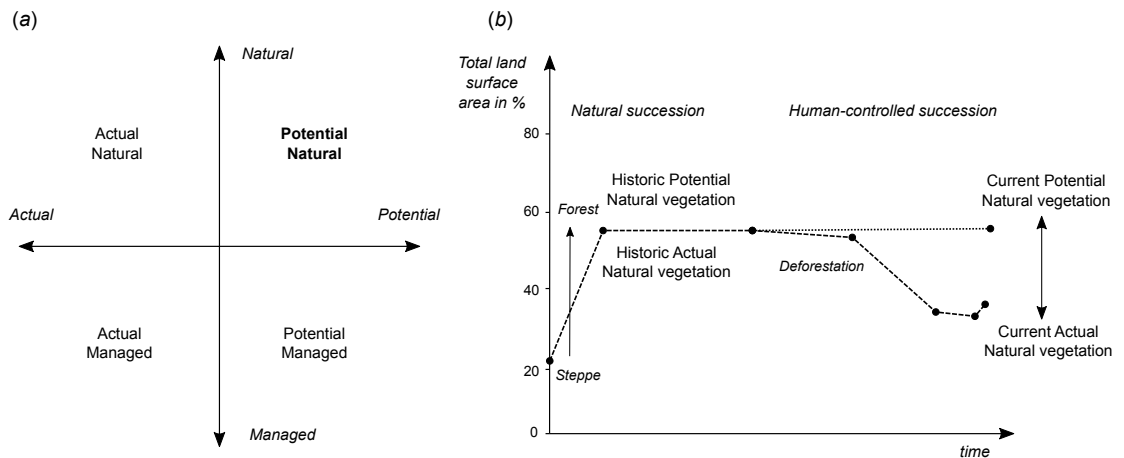


Figure 1. Schematic explanation of differences between (a) potential and actual natural/managed vegetation, and (b) current and historic vegetation in the context of global land area.

96 2. PNV model B: based on the autochthonous or native vegetation that includes also extinct species.

97 3. PNV model C: PNV based on any vegetation whether native or introduced or extinct.

98 Derivation of maps of PNV model A could be of interest to e.g. nature conservationists; PNV model
99 C could be of more interest to e.g. forestry and agroforestry organizations as it provides an objective basis
100 for introducing non-native species to a new area.

101 Conveniently, locations that have not been subject to human disturbance/management can provide
102 relevant information about vegetation cover in historic times, which can serve as a guide to PNV. A major
103 limitation of modeling PNV is that we unfortunately do not have equally detailed information about the
104 status of vegetation and environment across historic periods. For instance, about half of the Earth's mature
105 tropical forests have disappeared in the last 150 years and original habitats have been reduced to 10 %
106 (Hansen et al., 2013). Given that climates have changed and few areas are truly human impact "free",
107 even undisturbed historic vegetation only represents one possible expression of PNV for a given set of
108 climate conditions at a specific time.

109 Regardless of the hypothetical nature of PNV, the concept (both as a classification and as a regression
110 type problem) is still a helpful yardstick against which land cover change can be quantitatively measured
111 and land restoration designs can be planned. Erb et al. (2017) have estimated that almost half of the
112 standing global vegetation biomass carbon stocks has been lost, almost equally due to land cover change
113 (e.g. tree cover to cropland) and management effects within land cover types (e.g. croplands managed at
114 lower biomass carbon stocks than tree covered areas). PNV maps can thus help quantify such differences,
115 both deficit and surplus, in biomass stocks caused by the current land management system more objectively
116 and served as an input to the redesign of land management systems.

117 PNV mapping and species distribution modeling

118 In principle, PNV mapping is a special case of species distribution modeling (Elith and Leathwick, 2009;
119 Østbye Hemsing and Bryn, 2012; Hijmans and Elith, 2018): at the core of PNV mapping is statistical
120 modeling of the relationship between species (or natural associations of species or communities) and a

list of predictors i.e. biotic and abiotic site factors (Elith and Leathwick, 2009). The difference between mapping actual distribution of species and PNV mapping is that PNV involves extrapolating the model to the whole land mask, assuming a hypothetical distribution under a specific set of undisturbed bioclimatic and/or biophysical conditions:

$$\Pr(Y) = f(\text{Relief, BioClimate, Lithology}) \quad (1)$$

where Y is the target variable, which could be vegetation types or plant species with a finite number of states $Y \in \{1, 2, \dots, k\}$ and/or vegetation properties. PNV mapping can be considered as a *classification-type* or *regression-type* problem depending on whether we map factors such as vegetation types or continuous vegetation properties such as biomass or leaf area index.

The primary assumptions we make when applying a PNV model to the training data are:

1. The ecological gradients captured in training data reflect only natural ecological gradients and not human controls such as land use systems, civil engineering constructions, or one-off major disturbance events such as volcanic eruptions, floods, or tsunamis.
2. Remote sensing data such NDVI often reflect human-altered vegetation patterns and ought not be used as covariates in PNV mapping (Leong and Roderick, 2015).
3. The training data are representative of the study area, especially considering the feature space (ecological gradients) of the study area.

Assuming a log-linear relationship between ecological gradients and target variables, PNV classes can be modeled using a multinomial log-linear model:

$$f(k, i) = \beta_{0,k} + \beta_{1,k}x_{1,i} + \beta_{2,k}x_{2,i} + \dots + \beta_{M,k}x_{M,i} \quad (2)$$

where $f(k, i)$ is the linear predictor function, β are the regression coefficients associated with the m th explanatory variable and the k th outcome. An efficient implementation of the multinomial logistic regression is the `multinom` function from the R package `nnet` (Venables and Ripley, 2002). The output of predictions produced using `multinom` are k probability maps (0–100 %) such that all predictions at each site sum up to 1:

$$\sum_{k=1}^K \Pr(Y_i = k) = 1 \quad (3)$$

In this paper, all prediction models are used in the “probability” mode i.e. to derive probability maps per class.

Note that a PNV spatial prediction model divides geographic space among all possible states given the training points. It is therefore necessary, for Eq.(1), that all possible states of Y are represented with

148 training data so that the model can be applied over the whole spatial domain of interest. If all of the states
149 are not known, then the space will be artificially filled-in with known classes occupying similar ecological
150 niches and which can lead to prediction bias. In other words, as with species distribution modeling of
151 individual species, both presence and absence data play an equally important role for model calibration
152 (Elith and Leathwick, 2009).

153 **Input data: training points**

154 We consider three ground-truth data sets for model calibration:

- 155 1. an expanded version of the BIOME 6000 DB data set representing site-based reconstructions from
156 surface pollen samples of major vegetation types or biomes ([http://dx.doi.org/10.17864/
157 1947.99](http://dx.doi.org/10.17864/1947.99)),
- 158 2. EU Forest (Mauri et al., 2017) and GBIF (Global Biodiversity Information Facilities) occurrence
159 records of the 76 main forest tree taxa in Europe (<http://dx.doi.org/10.15468/dl.fhucwx>),
- 160 3. Long-term Fraction of Absorbed Photosynthetically Active Radiation (FAPAR) monthly images
161 derived using a time-series of Copernicus Global Land products ([https://land.copernicus.
162 eu](https://land.copernicus.eu)),

163 BIOME 6000 and EU Forest and GBIF occurrences are point data sets, while FAPAR consists of
164 remote sensing images at relatively fine spatial resolution (250 m), from which we sample a large number
165 of values (ca 100,000) using random sampling after masking for areas of natural vegetation.

166 **BIOME 6000**

167 The BIOME 6000 data set (<http://dx.doi.org/10.17864/1947.99>) includes vegetation reconstruc-
168 tions from modern pollen samples, preserved in lake and bog sediments and from moss polsters, soil and
169 other surface deposits. The use of pollen data to reconstruct PNV relies on the fact that although modern
170 pollen samples may contain markers of land use, the predominant pollen types found in any one sample
171 are those of the regional vegetation within a radius on the order of 10–30 km around the sampling site.
172 Even if forests have fragmented, these fragments continue to produce and disperse pollen grains, and the
173 composition of the pollen assemblage provides information on tree taxa that are still present.

174 The BIOME 6000 data set is an amalgamation of multiple data sets. BIOME 6000 initially produced
175 maps for individual regions: Europe, Africa and the Arabian Peninsula, the Former Soviet Union and
176 Mongolia and China. Additional regions were subsequently added including Beringia, western North
177 America, Canada and the eastern United States and Japan, and the data for northern Eurasia, China
178 and southern Europe and Africa were also updated. These regional compilations were summarized
179 in Prentice and Jolly (2000). Subsequent regional updates include China (Harrison et al., 2001), the
180 circum-Artic region (Bigelow et al., 2003), Australia (Pickett et al., 2004) and South America (Marchant
181 et al., 2009). Additionally, we have also combined these data with pollen-based vegetation reconstructions
182 from the Eastern Mediterranean-Black Sea-Caspian Corridor (EMBSecBIO) region (Marinova et al.,
183 2018) available from <http://dx.doi.org/10.17864/1947.109>, to produce a more complete and
184 up-to-date compilation of the BIOME 6000.

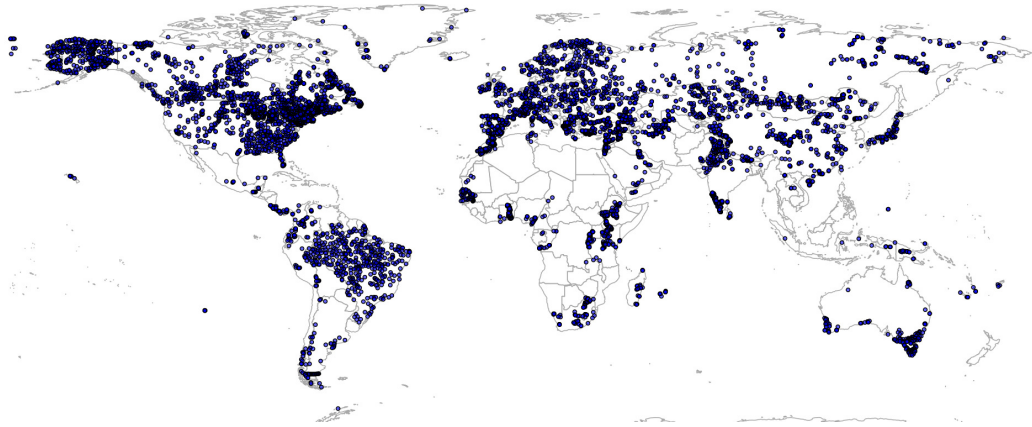


Figure 2. Spatial distribution of BIOME 6000 training points. A total of 8057 unique locations are shown on the map.

185 Some sites in the BIOME 6000 data set have multiple reconstructions based on multiple nearby
186 modern pollen samples (up to 30), which provides a useful measure of the reconstruction uncertainty, but
187 could lead to modeling bias because the number of modern samples varies between sites. To reduce these
188 unwanted effects, we use only the most frequently reconstructed biome at each site and for those sites
189 with two equally common reconstructions (ca. 900) we use both observations.

190 The number of biomes differentiated varies from region to region, and some biomes were only
191 reconstructed in specific regions where they are particularly characteristic, although they may occur, but
192 not be recognized, elsewhere. Furthermore, some biomes that can be recognized on the modern landscape
193 were never reconstructed in the BIOME 6000 data set (e.g. cushion forb tundra) — either because of
194 the sample distribution or because the characteristic plant-functional types were also spread amongst
195 other biomes. Simplified or “*megabiome*” classifications (e.g. [Harrison and Bartlein \(2012\)](#)) involve a
196 substantial loss of information. We have therefore created a new standardization of the classification
197 scheme (see further [Table 1](#); the final scheme has 20 globally applicable and distinctive biomes) which
198 preserve the maximum number of distinct biomes that were reconstructed as present in multiple regions.

199 There are relatively few data vegetation reconstructions for tropical South America, which could lead
200 to extrapolation problems and omission of important PNV classes in Latin America, but also potentially in
201 tropical parts of Africa and Asia. To reduce under-representation of tropics, we have added 350 randomly
202 simulated points based on the RADAM Brazil natural vegetation polygon map at high spatial detail
203 (Radam Vegetação SIRGAS map) ([Veloso et al., 1992](#)) obtained from <ftp://geofftp.ibge.gov.br/>.
204 Before generating the pseudo-observations for Brazil, we translated SIRGAS map legends to match the
205 BIOME 6000 classes. This translation is also available via the project’s github repository. This gave a
206 total of 8057 unique individual locations represented in the combined data set i.e. a total of 8959 training
207 observations ([Fig. 2](#)).

208 We have mapped the distribution of biomes for all land pixels, with the exception of water bodies,
209 barren land and permanent ice areas. Barren land and permanent ice areas were masked out using the
210 ESA’s global land cover maps for the period 2000–2015 (<https://www.esa-landcover-cci.org>)
211 and the long-term FAPAR images, both available at relatively fine resolution of 300 m. We only mask out

212 pixels that are permanent ice/barren ground and have a FAPAR = 0 throughout the period 2000–2015.

213 **European Forest Tree occurrence records**

214 For mapping PNV distribution of forest tree taxa (note: most of these are individual species, but some
 215 are only recognised at sub-genus or genus level) in Europe we have merged two point data sets: EU
 216 Forest (Mauri et al., 2017) (588,983 records covering 242 species) and GBIF occurrence records of
 217 the main forest tree taxa in Europe. The GBIF Occurrence data was downloaded on 23rd January
 218 2017 (<http://dx.doi.org/10.15468/dl.fhucwx>). We focus on modeling just the 76 forest tree taxa
 219 indicated in the European Atlas of Forest Tree Species (San-Miguel-Ayanz et al., 2016).

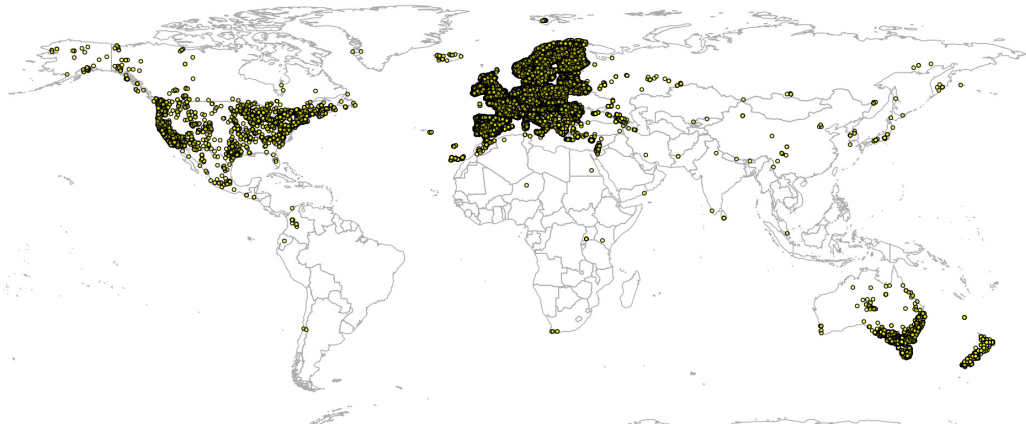


Figure 3. Merge of EU Forest (Mauri et al., 2017) and GBIF occurrence records used to build models to predict PNV for the 76 forest tree taxa. Total of 1,546,435 shown on the map.

220 Global GBIF occurrence data can be obtained by using the `rgbif` package, in which case the only
 221 important parameter is the `taxonKey` (e.g. “*Betula spp.*” corresponds to GBIF taxon key 2875008). After
 222 the bulk data download (which gives about 4 million occurrences), we imported all points and then subset
 223 occurrences based on the list of taxon keys and coordinate uncertainty (<2 km positional error). This
 224 gave a total of 1,546,435 training points from which about 2/3 are GBIF points (Fig. 3). We assume in
 225 further analysis that the EU Forest point locations and representativeness are more trustworthy, hence we
 226 assign 4× higher weights to these points than to the GBIF points.

227 Certain forest tree species (*Chamaecyparis lawsoniana*, *Eucalyptus globulus* and *Pseudotsuga men-*
 228 *ziesii*), that are shown in the European Atlas of Forest Tree Species are introduced i.e. planted and do not
 229 generally propagate naturally. Hence, they were removed from the list of target forest tree species. We
 230 retained, however, three species (*Ailanthus altissima*, *Picea sitchensis* and *Robinia pseudoacacia*) that are
 231 not native but are extensively naturalized. The total number of target forest tree taxa was 73.

232 We built predictive models for European forest tree taxa using information on their global distribution,
 233 but only generate predictions for Europe. In other words, we use a global compilation for model training
 234 to increase the precision of the definition of the ecological niche of each taxon, but then predict only for
 235 Europe as the selection of taxa is based on the European Atlas of Forest Tree Species (San-Miguel-Ayanz
 236 et al., 2016).

237 **FAPAR**

238 Fraction of Absorbed Photosynthetically Active Radiation (FAPAR) monthly images for 2014–2017 were
 239 obtained from <https://land.copernicus.eu> (original values reported in the range 0–235 with scaling
 240 factor 1/255 i.e. with a maximum value of 0.94). From a total of 142 images downloaded from <https://land.copernicus.eu> we derived minimum, median and maximum value of FAPAR per month (12)
 241 using the 95 % probability interval using the data.table package (<http://r-datatable.com>). For
 242 regression modeling we only report results of predictions of median values of FAPAR; predictions of
 243 minimum and maximum FAPAR can be obtained from the data repository.
 244

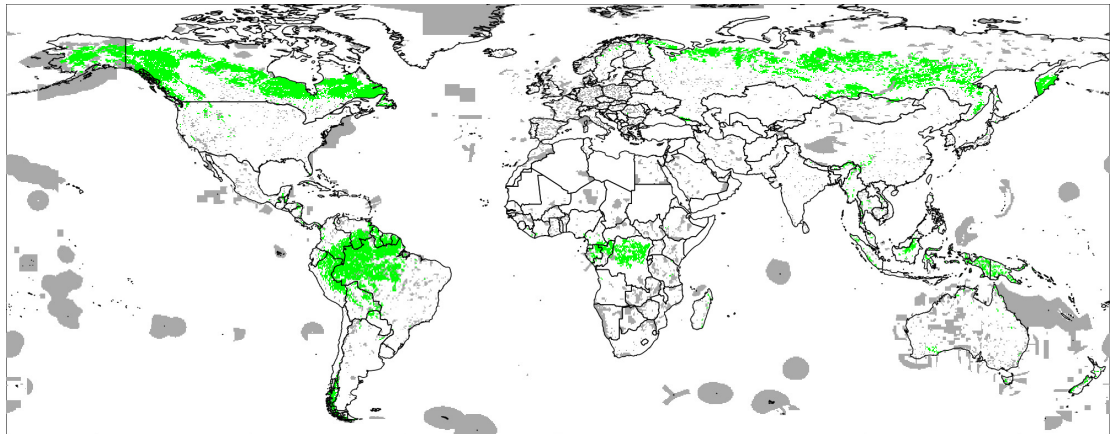


Figure 4. World's Protected Areas (dark gray) based on <http://protectedplanet.net> and Intact Forest Landscapes for year 2000 (green) based on <http://intactforests.org>. These maps were used to randomly select some 30,000 training points to predict potential FAPAR under PNV.

245 We model median and upper 95 % FAPAR values as a function of the same covariate layers used in
 246 all three case studies. For model training we use ca. 30,000 randomly sampled points (Simple Random
 247 Sampling) exclusively from protected area as shown in the World Database on Protected Areas (WDPA)
 248 data set (<http://protectedplanet.net>) and the Intact Forest Landscapes (IFL) data set for 2000 and
 249 2013 (Potapov et al., 2008) Fig. 4). We use about $3\times$ more training points from the IFL 2013 areas for
 250 model development than from the WDPA and IFL 2000 masks to emphasize more ecological conditions
 251 of intact vegetation.

252 The prediction model for FAPAR under PNV is in the form of:

$$R> \text{FAPAR} \sim c_m + X_{1m} + X_{2m} + X_3 + \dots + X_p$$

253 where X_{1m} is the covariate with monthly values (for example precipitation, day-time and night-time
 254 temperatures etc), X_3 is the environmental covariates that do not vary through year (e.g. lithology or DEM
 255 derivatives), and c_m is the cosine of the month number:

$$c_m = \cos(\mu/12 \cdot 2 \cdot \pi) \quad (4)$$

256 where μ is the month number 1–12. The total number of training observations used to build models is in
 257 fact 180,483 (each training site is represented up to 12 times).

258 For PNV FAPAR mapping we have masked out all water bodies including lakes and rivers, following
259 the ESA's global land cover maps for the period 2000–2015 (<https://www.esa-landcover-cci.org>)
260 and permanent ice/barren ground.

261 **Input data: environmental covariates**

262 For modeling purposes, we use a stack of 160 spatially explicit co-variate data layers that represent
263 standard ecological gradients essential for growth and survival of plants:

- 264 • DEM derivatives quantifying various landscape metrics and hydrological processes: slope, curva-
265 ture, topographic index, topographic openness, valley depth and multi-resolution valley bottom
266 index; all derived using the SAGA GIS (Conrad et al., 2015);
- 267 • Mean, minimum and maximum monthly temperatures derived as a mean between WorldClim v2
268 (<http://worldclim.org/version2>) and CHELSA climate (Karger et al., 2017).
- 269 • Mean monthly precipitation images derived as a weighted average between the WorldClim v2,
270 CHELSA climate and Global Precipitation Measurement Integrated Multi-satellitE Retrievals for
271 GPM (IMERG) rainfall product.
- 272 • CHELSA Bioclimatic layers downloaded from <http://chelsa-climate.org/>, including: an-
273 nual mean temperature, mean diurnal temperature range, isothermality (day-to-night temperature
274 oscillations relative to the summer-to-winter oscillations), temperature seasonality (standard de-
275 viation of monthly temperature averages), maximum temperature of warmest month, minimum
276 temperature of coldest month, temperature annual range, mean temperature of warmest quarter,
277 mean temperature of coldest quarter, annual precipitation amount, precipitation of wettest month,
278 precipitation of driest month, precipitation of wettest quarter, precipitation of driest quarter (Karger
279 et al., 2017);
- 280 • European Space Agency's CCI-LC snow probability monthly averages based on MODIS snow
281 products MOD10A2 downloaded from [http://maps.elie.ucl.ac.be/CCI/viewer/index.
282 php](http://maps.elie.ucl.ac.be/CCI/viewer/index.php);
- 283 • USGS Global Ecophysiology landform classification and lithological map at 250 m resolution
284 obtained from <http://rmgsc.cr.usgs.gov/outgoing/ecosystems/Global/> and based on
285 Global Lithological Map (GLiM) (Hartmann and Moosdorf, 2012);
- 286 • MODIS Cloud fraction monthly images obtained from <http://www.earthenv.org/cloud> (Wil-
287 son and Jetz, 2016);
- 288 • Global Water Table Depth in meters based on Fan et al. (2013);
- 289 • NASA's monthly MODIS Precipitable Water Vapor images (MYDAL2_M_SKY_WV data set at <http://neo.sci.gsfc.nasa.gov>);
- 290 • Potential wetlands GIEMS map (Fluet-Chouinard et al., 2015);
- 291 • Global Surface Water dynamics images: occurrence probability, surface water change, and
292 water maximum extent; downloaded from [https://global-surface-water.appspot.com/
293 download](https://global-surface-water.appspot.com/download) (Pekel et al., 2016);
- 294 • Density of earthquakes based on the USGS Earthquake Archives ([http://earthquake.usgs.
295 gov/earthquakes/](http://earthquake.usgs.gov/earthquakes/));
- 296

297 Some CHELSA bioclimatic layers contained too many missing pixels or artifacts (e.g. mean temper-
298 ature of wettest quarter, mean temperature of driest quarter, precipitation seasonality, precipitation of
299 warmest quarter and precipitation of coldest quarter) and hence were not used for further modeling to
300 avoid propagating those artifacts to final predictions.

301 All original layers have been resampled to the standard grid at a spatial resolution of 1/120 decimal
302 degrees (about 1 km) covering latitudes between -62.0 and 87.37. Some layers such as Water Vapour
303 needed to be downscaled from 10 km to 1 km resolution, for which we used the bicubic splines algorithm
304 as implemented in GDAL (Mitchell and GDAL Developers, 2014). We do not map Antarctica as this
305 continent is dominantly covered with permanent ice and there are no training points. We limit all analysis
306 to 1 km i.e. 1/120 degrees in geographical coordinates, to avoid too high of a computational load, even
307 though many of environmental covariates are also available at finer resolutions.

308 We use the same stack of covariates for mapping global distribution of biomes, FAPAR and forest tree
309 species in Europe, in order to be able to compare model performance and investigate whether the most
310 important covariates differ among the three case studies.

311 Machine Learning Algorithms (MLA) examined

312 We examine predictive performance of the following MLA's:

- 313 • Neural networks (Venables and Ripley, 2002),
- 314 • Random forest (Breiman, 2001; Cutler et al., 2007; Biau and Scornet, 2016; Hengl et al., 2018),
- 315 • Generalized Boosted Regression Models (Friedman, 2002),
- 316 • K-nearest neighbors (Venables and Ripley, 2002),

317 Neural networks are available from several packages in R. Here we use the nnet package (Ripley and
318 Venables, 2017) also described in Venables and Ripley (2002). Random forest is efficiently implemented in
319 the ranger package (Wright and Ziegler, 2016) and can be used to process large data sets. Generalized
320 Boosted Regression Models are available via the gbm package (Ridgeway, 2017). The K-nearest
321 Neighbour Regression is available via the class package i.e. the knn function (Venables and Ripley,
322 2002). Of these four algorithms, the K-nearest neighbors is computationally the least intensive and results
323 in relatively simple models, while random forest is computationally the most intensive and results in
324 large models. However, a limitation of the K-nearest neighbors approach is that it does not handle high
325 dimensional data in comparison to random forest or neural nets.

326 We also test using the same packages to fit models for regression-type problems (e.g. modeling
327 of FAPAR), with the exception of the class package i.e. the knn function which can only be used for
328 classification problems. For modeling FAPAR we instead added use of the Cubist approach, available via
329 the Cubist package (Kuhn et al., 2017), and the Extreme Gradient Boosting approach available via the
330 xgboost package (Chen and Guestrin, 2016).

331 The caret package has many more MLA of interest for classification and regression problems than
332 presented here, but many are not fully optimized for large data sets and hence also not applicable for large
333 data sets (\gg 1000 observations with \gg 100 covariates).

334 **Model selection**

335 For model fitting and model selection we use the caret package implementation for automated evaluation
336 of models. When comparing performance of the models we look at classification accuracy based on
337 cross-validation with refitting implemented in the caret package via the setting (Kuhn, 2008; Kuhn and
338 Johnson, 2013):

```
R> ctrl <- trainControl(method="repeatedcv", number=5, repeats=2)
```

339 which translates as: models are refit 5 times using 80 % of the data and predictions derived from the fitted
340 models are compared with the remaining observations; this process is then repeated two times to produce
341 stable results. The reported accuracy is the map accuracy (0–100 %) and/or Root Mean Square Error
342 (RMSE) derived using all merged cross-validations (Kuhn, 2008; Kuhn and Johnson, 2013). Since most
343 of the data sets are fairly large and model fitting can take hours, even in a High Performance Computing
344 environment, we limit the number of repetitions to 2.

345 For FAPAR (regression modeling) and selection of the final prediction model we use the same repeated
346 cross-validation as implemented via the caret package. This is, in principle, similar to evaluation of the
347 classification accuracy, except the comparison criterion is RMSE.

348 All analyses were run on a High Performance Computing Amazon ec2 server with 64 threads (32
349 CPU's) and 256 GiB RAM. Total computing time to produce all outputs is about 12 hours of optimized
350 computing (or about 600 CPU hours). 1 km data can be processed with 2 degree tiles, which usually
351 requires some 5000 tiles to represent the land mask. All processing steps and preparation of input
352 and output maps are fully documented at <https://github.com/envirometrix/PNVmaps>. All output
353 maps are available for download via <http://dx.doi.org/10.7910/DVN/QQHCIK> under the Open
354 Database License (ODbL).

355 **Performance of classification algorithms**

356 Performance of classification algorithms is assessed using 5-fold cross-validation with refitting of models.
357 For evaluation of the mapping accuracy for biomes and tree species we use the map purity (0–100 %) and
358 kappa metrics for the dominant (hard) classes as the key measures of predictive performance (Kuhn
359 and Johnson, 2013). For each class we also provide predicted probabilities, which can be used to model
360 transition zones and correlation between classes. For the predicted probabilities of class occurrences (0–1)
361 we derived the True Positive Rate (TPR) and the Area Under the receiver operating characteristic Curve
362 (AUC) as implemented in the ROCR package (Sing et al., 2005, 2016). TPR value = 1 indicates a perfect
363 match to the class positives in ground data while TPR values < 0.5 can be considered poor mapping
364 accuracy. Likewise, values of AUC close to 1 indicate high prediction performance, while values around
365 0.5 and below are considered poor. TPR and AUC provide probably a more informative measure of the
366 mapping accuracy than overall mapping accuracy / kappa, as they also allow detection of problematic
367 classes.

368 We also use Scaled Shannon Entropy Index, which can be derived using the per-class probability maps

369 (Shannon, 1949; Borda, 2011):

$$\text{SSEI}_s(x) = - \sum_{i=1}^b P_i(x) \cdot \log_b P_i(x) = \frac{- \sum_{i=1}^b P_i(x) \cdot \log P_i(x)}{-b \cdot b^{-1} \cdot \log b^{-1}} \quad (5)$$

370 where b is the total number of possible classes and P is probability of class i . The Scaled Shannon
 371 Entropy Index (SSEI) is in the range from 0–1, where 0 indicates a perfect classification and 1 (or 100 %)
 372 indicates maximum confusion. Scaled Shannon Entropy Index should not be confused with classification
 373 accuracy assessment. For example, $\text{SSEI}_s < 60\%$ indicates relatively low confusion between classes i.e.
 374 high accuracy, while mapping error of 60 % would be considered a relatively poor classification accuracy
 375 result.

376 For the biomes data set, where spatial clustering of points is significant, we also use repeated spatial
 377 cross-validation as implemented in the mlr package (Bischl et al., 2016):

```
R> learner.rf = makeLearner("classif.ranger", predict.type = "prob")
R> resampling = makeResampleDesc("SpRepCV", fold = 5, reps = 5)
```

378 It has been shown that spatial autocorrelation in data and serious spatial clustering in training points
 379 can lead to somewhat biased estimate of the actual accuracy (Brenning, 2012). A solution to this problem
 380 is to apply spatial partitioning so that possible bias due to spatial proximity is minimized.

381 We also compare results of modeling potential distribution of tree species in Europe with the habitat
 382 type maps of Europe produced independently by San-Miguel-Ayanz et al. (2016) and Brus et al. (2012).
 383 This comparison is visually based only.

384 Performance of regression algorithms

385 Performance of regression algorithms is also assessed using 5–fold cross-validation with refitting of
 386 models. For assessment of the mapping accuracy for FAPAR we use as the main performance measures
 387 the root mean squared error (RMSE):

$$\text{RMSE} = \sqrt{\frac{\sum_{j=1}^m [\hat{y}(\mathbf{s}_j) - y(\mathbf{s}_j)]^2}{n}} \quad (6)$$

388 and mean error (ME):

$$\text{ME} = \frac{\sum_{j=1}^m [\hat{y}(\mathbf{s}_j) - y(\mathbf{s}_j)]}{n} \quad (7)$$

389 where $\hat{y}(\mathbf{s}_j)$ is the predicted value of y at the cross-validation location, and m is total number of cross-
 390 validation points. We also report amount of variation explained by the model (R^2) derived as:

$$R^2 = \left[1 - \frac{\text{SSE}}{\text{SST}} \right] \times 100\% \quad (8)$$

391 where SSE is the sum of squared errors at cross-validation points and SST is the total sum of squares. A
 392 coefficient of determination close to 1 indicates a perfect model.

393 RESULTS

394 Global maps of biomes

395 Results showed that a relatively accurate model of PNV could be produced from the BIOME 6000 data
 396 set using the existing stack of covariates at 1 km spatial resolution. Results of cross-validation show
 397 the random forest (RF) model to be the best performing method and distinctively superior to all other
 398 approaches (Fig. 5). The choice of the random forest $mtry$ parameter had little impact on overall accuracy,
 399 most likely because there was a high overlap in covariate maps so that even with smaller $mtry$ bagging
 400 the performance was relatively similar. The best prediction accuracy from among the four methods used
 401 for mapping global biomes was about 68%. The predicted biome classes are presented in Fig. 6.

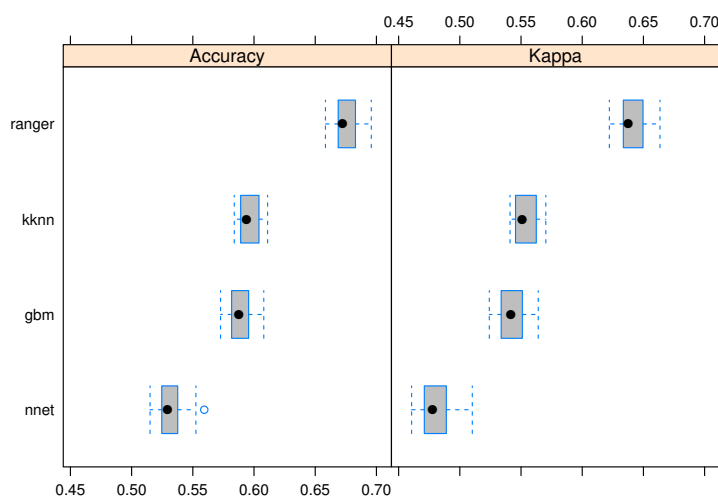


Figure 5. Predictive performance of the target machine learning algorithms for mapping global distribution of biomes ($N = 8653$; spatial distribution of training points is available in Fig. 2). ranger = random forest, kkn = K-nearest neighbors, gbm = Generalized Boosted Regression Models, nnet = Neural networks.

402 The most important covariates for the random forest model were: total annual precipitation, monthly
 403 temperatures, CHELSA bioclimatic layers, atmospheric water vapor images and monthly precipitation.
 404 Landform parameters and lithology are not amongst the top 20 most important predictors. The decline in
 405 variable importance was, however, gradual — even lower ranked covariates might still affect the accuracy
 406 of predictions.

407 The detailed cross-validation results show that the only difficult class to predict was prostrate dwarf
 408 shrub tundra (Table 1). The TPR value for most class probabilities ranges from 0.83 to 0.94 indicating
 409 relatively high match with ground data. The Scaled Shannon Entropy Index map (Fig. 7) showed that the
 410 zones of highest confusion between classes can be found in Afghanistan, Nepal, mountainous parts of the
 411 USA and Mexico, parts of Angola and Zambia. The map of the SSEI is comparable to the confusion map

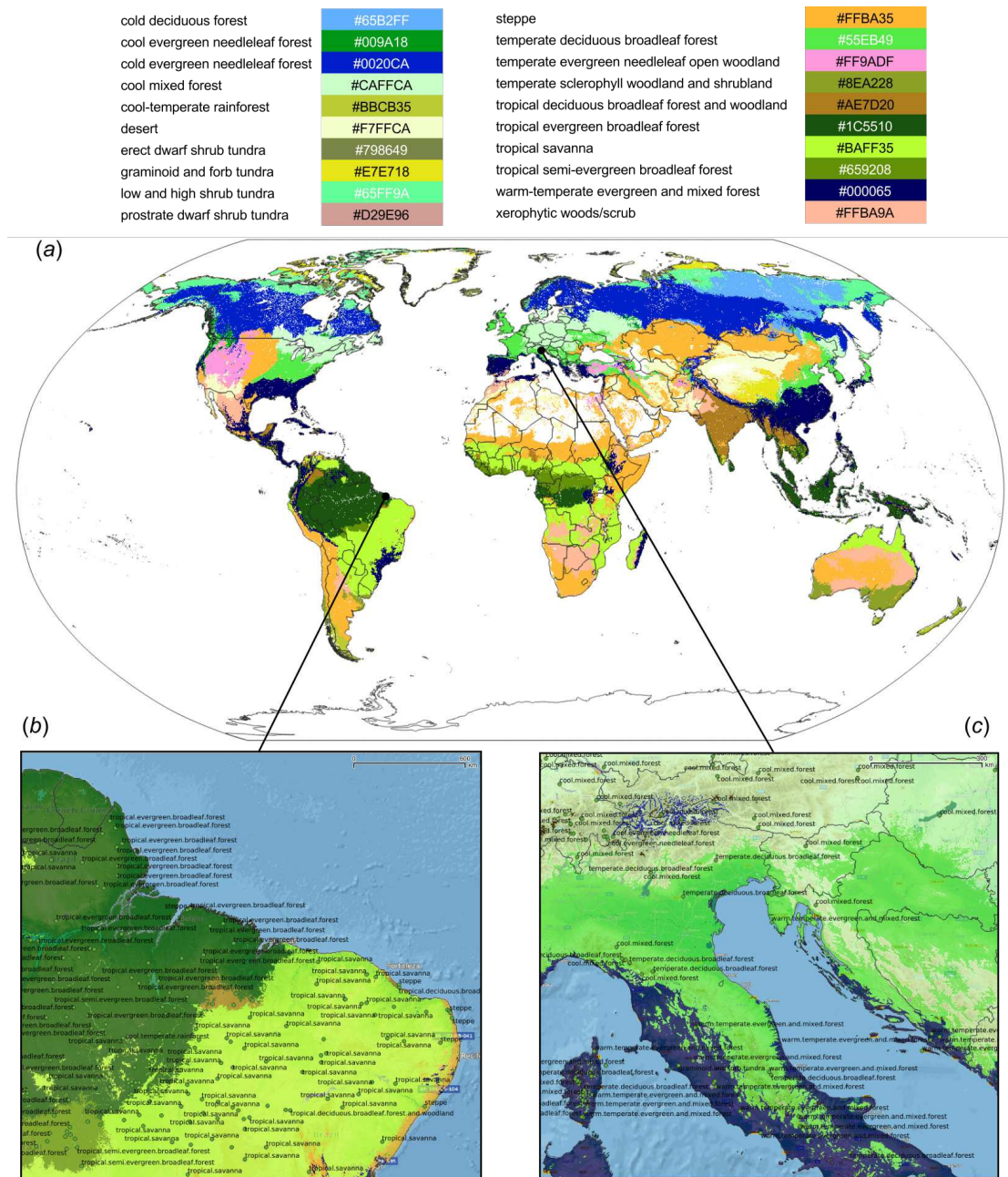


Figure 6. Predicted PNV distribution for (a) global biomes with a zoom in on areas in Brazil (b) and Europe (c). Labels indicates training points from the BIOME 6000 data set (Fig. 2). Background map data: Google, DigitalGlobe.

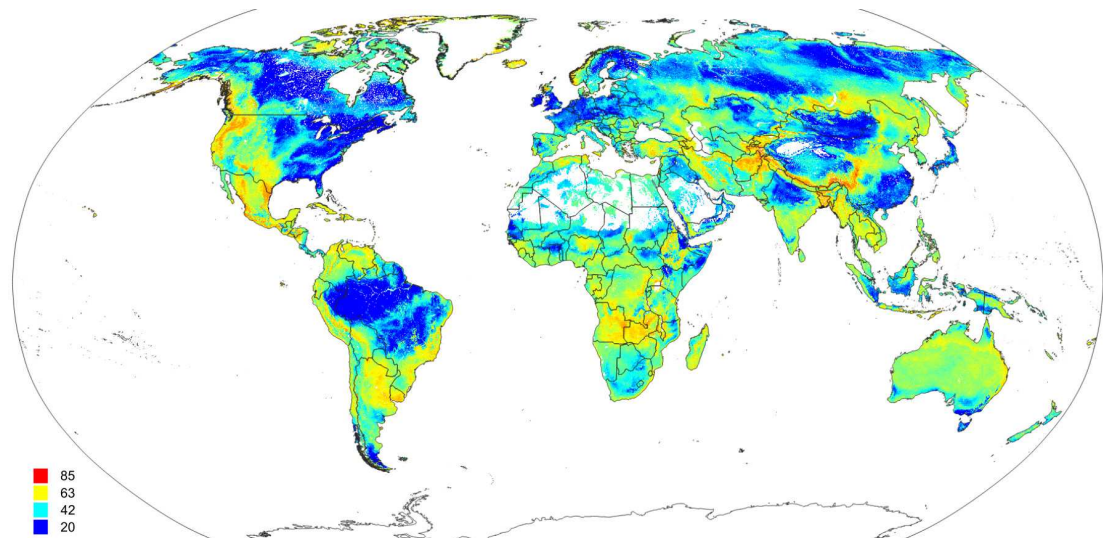


Figure 7. Scaled Shannon Entropy Index (SSEI) derived using predicted probabilities for 20 biomes (classes) based on Eq.(5). High values of SSEI (red color) indicate high confusion between classes.

Table 1. Summary results of cross-validation for mapping global distribution of biomes (20 classes). Classification accuracy for predicted class probabilities is based on 5-fold cross-validation with refitting. ME = “Mean Error”, TPR = “True Positive Rate”, AUC = “Area Under Curve”, N = “Number of occurrences”.

Biome class	ME	TPR	AUC	N
cold deciduous forest	-0.01	0.89	0.96	201
cold evergreen needleleaf forest	0.01	0.87	0.98	892
cool evergreen needleleaf forest	-0.07	0.87	0.93	201
cool mixed forest	0.01	0.86	0.97	1549
cool temperate rainforest	0.01	0.92	0.99	95
desert	0.00	0.89	0.96	330
erect dwarf shrub tundra	-0.01	0.89	0.98	145
graminoid and forb tundra	-0.03	0.83	0.91	128
low and high shrub tundra	-0.01	0.88	0.98	393
prostrate dwarf shrub tundra	-0.02	0.54	0.90	11
steppe	0.01	0.87	0.94	889
temperate deciduous broadleaf forest	-0.01	0.84	0.94	961
temperate evergreen needleleaf open woodland	0.01	0.92	0.97	307
temperate sclerophyll woodland and shrubland	0.00	0.94	0.99	154
tropical deciduous broadleaf forest and woodland	0.01	0.86	0.97	215
tropical evergreen broadleaf forest	0.00	0.87	0.99	333
tropical savanna	0.01	0.89	0.99	291
tropical semi evergreen broadleaf forest	-0.05	0.87	0.98	160
warm temperate evergreen and mixed forest	0.01	0.85	0.96	985
xerophytic woods scrub	-0.02	0.88	0.95	388

412 produced by [Levvasseur et al. \(2012\)](#), except in our case the Rocky Mountains in USA and mountains
 413 chains in South America show somewhat higher confusion. Many of the areas with high confusion index
 414 occur because the prediction model has problems distinguishing between closely-related biomes such as
 415 the “cold evergreen needleleaf forest” and “cool evergreen needleleaf forest” (e.g. Scotland).

416 Results of the accuracy assessment based on the spatial Cross-Validation (mlr package implementation

417 (Bischi et al., 2016)) further indicate that the spatial clustering of points does have a large effect on the
 418 mapping accuracy: spatial CV drops from 0.68 to 0.33 and weighted kappa to 0.45. This likely happens
 419 due to high spatial clustering of the biome points and due to the high spatial autocorrelation of biomes.

420 European forest tree species

421 The results of 5-fold cross validation with re-fitting at each fold, confirms that random forest was also the
 422 best prediction method for the forest taxa data set (Fig. 8). The overall mapping accuracy was significantly
 423 lower than for biomes, but this reduction in accuracy was to be expected as many of these taxa occur
 424 in communities, resulting in natural overlap of forest tree taxa distribution. The mapping accuracy of
 425 individual taxa, however, can be relatively high with TPR values of between 0.16–0.90 and an average
 426 value of around 0.69 (Table 2). The final maps (Fig. 9) showed a relatively good match with ground
 427 data, meaning that with the exception of some species of rarer occurrence (*Picea omorika*, *Cupressus*
 428 *sempervirens*, *Prunus mahaleb*), the species probability distribution maps were relatively accurate.

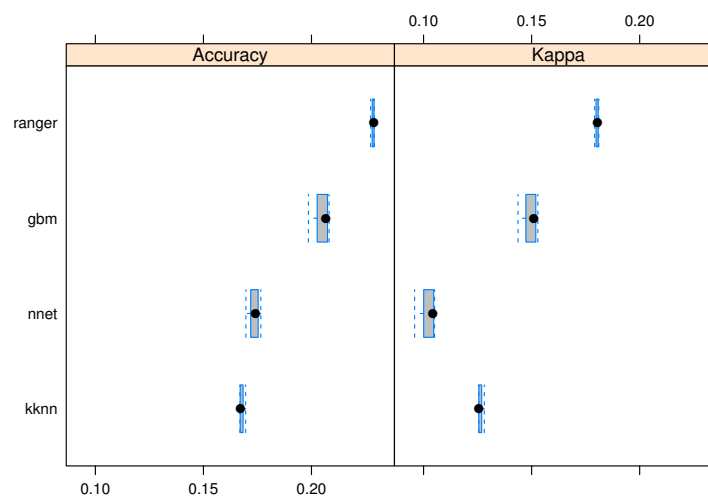


Figure 8. Predictive performance of the target machine learning algorithms for mapping forest tree species ($N = 1.5$ million distribution of training points is available in Fig. 3). ranger = random forest, gbm = Generalized Boosted Regression Models, nnet = Neural networks, kkn = K-nearest neighbors.

Table 2. Results of cross-validation for the forest tree taxa. Classification accuracy for predicted class probabilities based on 5-fold cross-validation. ME = “Mean Error”, TPR = “True Positive Rate”, AUC = “Area Under Curve”, N = “Number of occurrences”. Taxa with less than < 50 observations were omitted from analysis.

Species name	GBIF taxon ID	ME	TPR	AUC	N
<i>Abies alba</i>	2685484	-0.01	0.77	0.92	16,150
<i>Acer campestre</i>	3189863	-0.01	0.65	0.83	19,819
<i>Acer platanoides</i>	3189846	-0.02	0.68	0.82	30,801
<i>Acer pseudoplatanus</i>	3189870	-0.01	0.69	0.79	65,039
<i>Aesculus hippocastanum</i>	3189815	-0.01	0.59	0.85	8,088
<i>Ailanthus altissima</i>	3190653	0.04	0.69	0.92	1,576
<i>Alnus cordata</i>	2876607	0.05	0.73	0.95	904
<i>Alnus glutinosa</i>	2876213	0.00	0.71	0.77	91,292

Continued on next page

Table 2 – Continued from previous page

Species name	GBIF taxon ID	ME	TPR	AUC	N
<i>Alnus incana</i>	2876388	-0.03	0.76	0.95	6,873
<i>Betula spp.</i>	2875008	-0.03	0.63	0.83	7,313
<i>Carpinus betulus</i>	2875818	0.00	0.75	0.89	22,765
<i>Carpinus orientalis</i>	2875780	0.07	0.21	0.92	284
<i>Castanea sativa</i>	5333294	0.00	0.74	0.91	13,049
<i>Celtis australis</i>	2984492	-0.01	0.54	0.92	594
<i>Cornus mas</i>	3082263	0.03	0.51	0.90	827
<i>Cornus sanguinea</i>	3082234	-0.03	0.59	0.82	8,837
<i>Corylus avellana</i>	2875979	-0.02	0.67	0.76	48,140
<i>Cupressus sempervirens</i>	2684030	-0.04	0.21	0.70	284
<i>Euonymus europaeus</i>	3169131	-0.02	0.61	0.83	12,119
<i>Fagus sylvatica</i>	2882316	0.00	0.73	0.81	89,044
<i>Frangula alnus</i>	3039454	-0.02	0.71	0.86	26,873
<i>Fraxinus angustifolia</i>	7325877	-0.05	0.63	0.94	1,757
<i>Fraxinus excelsior</i>	3172358	0.00	0.67	0.74	91,111
<i>Fraxinus ornus</i>	3172347	0.02	0.86	0.99	2,765
<i>Ilex aquifolium</i>	5414222	-0.01	0.66	0.82	26,873
<i>Juglans regia</i>	3054368	-0.03	0.60	0.89	3,643
<i>Juniperus communis</i>	2684709	-0.03	0.71	0.86	21,189
<i>Juniperus oxycedrus</i>	2684451	-0.07	0.71	0.97	1,705
<i>Juniperus phoenicea</i>	2684640	-0.07	0.74	0.98	1,137
<i>Juniperus thurifera</i>	2684528	-0.03	0.87	0.99	1,886
<i>Larix decidua</i>	2686212	-0.01	0.71	0.89	15,581
<i>Olea europaea</i>	5415040	0.00	0.90	0.99	7,080
<i>Ostrya carpinifolia</i>	5332305	0.06	0.90	0.99	1,809
<i>Picea abies</i>	5284884	0.02	0.76	0.86	122,713
<i>Picea sitchensis</i>	5284827	0.05	0.80	0.96	13,023
<i>Pinus cembra</i>	5285134	-0.01	0.77	0.96	853
<i>Pinus halepensis and Pinus brutia</i>	5285604	0.03	0.86	0.99	16,951
<i>Pinus mugo</i>	5285385	0.00	0.85	0.98	6,667
<i>Pinus nigra</i>	5284809	0.01	0.79	0.93	13,540
<i>Pinus pinaster</i>	5285565	0.01	0.86	0.98	17,080
<i>Pinus pinea</i>	5285165	-0.04	0.85	0.99	4,910
<i>Pinus sylvestris</i>	5285637	0.02	0.78	0.85	153,928
<i>Populus alba</i>	3040233	-0.01	0.54	0.86	4,522
<i>Populus nigra</i>	3040227	-0.01	0.65	0.89	5,478
<i>Populus tremula</i>	3040249	-0.02	0.66	0.74	44,057
<i>Prunus avium</i>	3020791	-0.01	0.63	0.77	25,711
<i>Prunus cerasifera</i>	3021730	0.00	0.73	0.94	3,928
<i>Prunus mahaleb</i>	3022789	-0.01	0.31	0.75	517
<i>Prunus padus</i>	3021037	-0.03	0.63	0.78	21,705
<i>Prunus spinosa</i>	3023221	-0.01	0.69	0.81	31,783
<i>Quercus cerris</i>	2880580	0.00	0.80	0.97	4,109
<i>Quercus ilex</i>	2879098	0.02	0.85	0.99	22,972
<i>Quercus pubescens</i>	2881283	0.01	0.86	0.98	9,096
<i>Quercus pyrenaica</i>	2878826	0.00	0.88	0.99	6,253
<i>Quercus robur and Quercus petraea</i>	2878688	0.01	0.69	0.76	141,938
<i>Quercus suber</i>	2879411	-0.04	0.86	0.99	5,504
<i>Robinia pseudoacacia</i>	5352251	0.01	0.71	0.90	13,411
<i>Salix alba</i>	5372513	0.02	0.72	0.90	11,938
<i>Salix caprea</i>	5372952	-0.03	0.68	0.78	40,879
<i>Sambucus nigra</i>	2888728	0.00	0.70	0.81	44,961
<i>Sorbus aria</i>	3012680	-0.01	0.59	0.87	5,426
<i>Sorbus aucuparia</i>	3012167	-0.01	0.70	0.76	86,977
<i>Sorbus domestica</i>	3013206	-0.04	0.48	0.87	801
<i>Sorbus torminalis</i>	3012567	-0.03	0.62	0.92	2,558
<i>Taxus baccata</i>	5284517	-0.02	0.58	0.82	8,062
<i>Tilia spp.</i>	3152041	-0.02	0.50	0.82	4,393
<i>Ulmus spp.</i>	2984510	-0.03	0.64	0.92	5,426
<i>Tilia spp.</i>	3152041	0.00	0.58	0.85	4,522
<i>Ulmus spp.</i>	2984510	-0.02	0.69	0.91	5,375

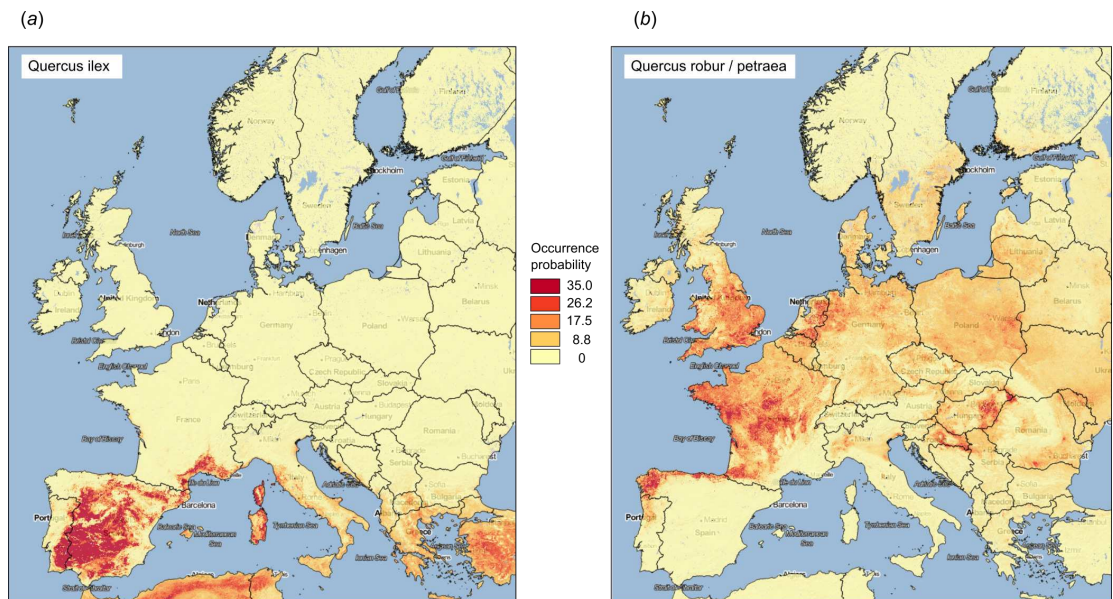


Figure 9. Examples of predicted PNV distributions (probabilities) for European forest tree species (a) *Quercus Ilex* (GBIF ID: 2879098; 36,724 training points) and (b) *Quercus robur / petraea* (GBIF ID: 2878688; 404,296 training points). Background map data: Google, DigitalGlobe.

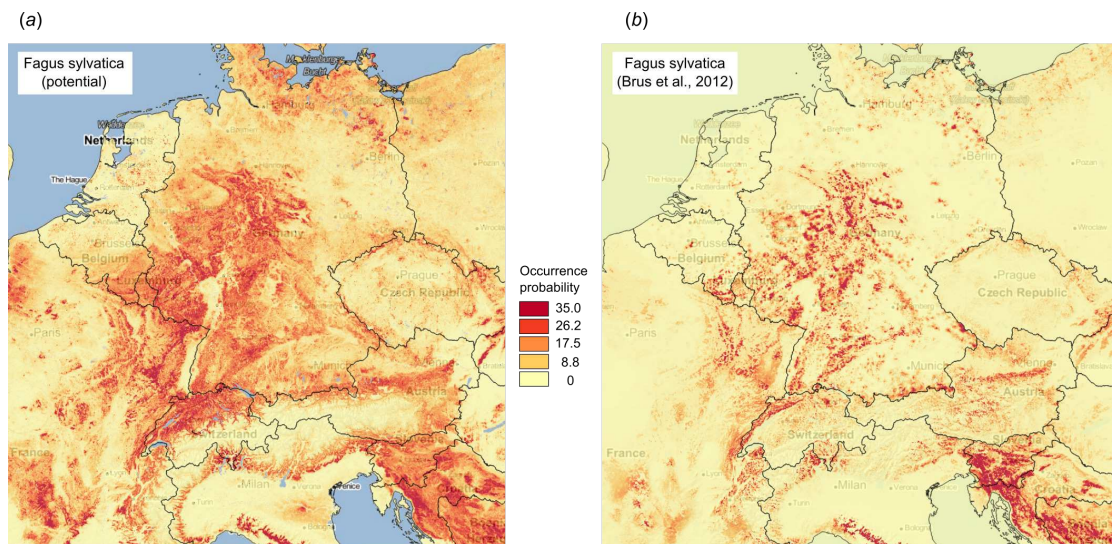


Figure 10. Comparison between predicted PNV distribution for (a) *Fagus sylvatica* (GBIF ID: 2882316) based on our results, and (b) based on the maps generated by Brus et al. (2012) i.e. showing the presumed actual distribution of the tree species. Background map data: Google, DigitalGlobe.

429

430 The most important predictors in the random forest model for forest tree taxa were mean annual daily
 431 temperature, other monthly temperatures, elevation, CHELSA bioclimatic images, monthly precipitation
 432 and MODIS cloud fraction images. Covariates for lithology and landform classification did not feature in
 433 the top 20 predictors. It could be that the Global Lithological Map (GLiM) (Hartmann and Moosdorf,
 434 2012), which was used to represent changes in lithology, is too general for this scale of work.

435 Fig. 10 illustrates differences between the map of actual distribution of *Fagus sylvatica*, generated by

436 Brus et al. (2012), and our predictions. In this case, the potential for extending habitat of *Fagus sylvatica*
437 is significant, especially over parts of France and Germany.

438 Correlation analysis using all predicted distribution maps (matrix of Pearson's rho rank correlation
439 coefficients for all possible pairs) indicated that many forest species are positively correlated, especially
440 *Fagus sylvatica* and *Abies alba* and *Populus nigra* and *Salix alba*. High overlap between species probability
441 maps reflects co-existence within communities, and thus could help with objectively defining forest
442 communities.

443 **Global monthly FAPAR**

444 The random forest approach also produced the best predictions of potential FAPAR (Fig. 11). The models
 445 for FAPAR were highly significant with R-squared around 90 % and RMSE at ± 24 (original values in
 446 the range 0–232 where 235 corresponds to FAPAR=100 %) for the most accurate model based on 5–fold
 447 Leave-Location-Out cross-validation. However, unlike with biomes and forest species distributions, the
 448 regression-tree Cubist model achieves equal performance to that of random forest. The most important
 449 covariates for predicting FAPAR were total annual precipitation, MODIS cloud fraction images, CHELSA
 450 bioclimatic images, and monthly precipitation images. The caret package further suggested that `mtry`
 451 parameter for Random Forest needs to be set higher than the default values for modeling FAPAR. Setting
 452 up `mtry` >25 helps reduce the RMSE by about 7–8 %.

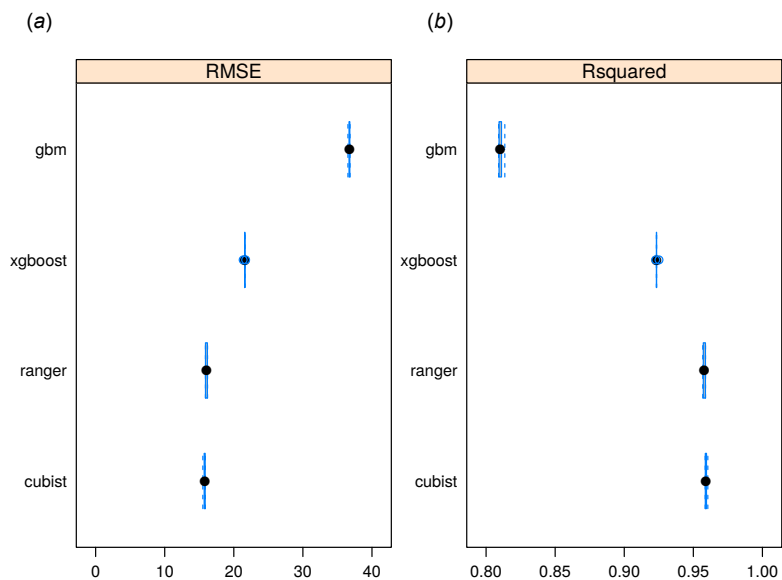


Figure 11. Predictive performance of four machine learning algorithms for mapping global distribution of FAPAR ($N = 180,990$). gbm = Generalized Boosted Regression Models, xgboost = Extreme Gradient Boosting, ranger = random forest, cubist = Cubist Regression Models. (a) RMSE = Root Mean Square Error, (b) R-squared.

453 Fig. 12 depicts an example of actual vs predicted (PNV based) FAPAR for February in the urban
 454 area around São Paulo, where lower actual FAPAR reflects the removal of natural vegetation. Even
 455 larger differences between the potential and actual FAPAR are observed in parts of Africa (Fig. 13),
 456 likely reflecting land degradation and destruction of vegetation cover. In areas of intensive agricultural
 457 production (e.g. Western Australia and Midwest USA), actual FAPAR can be much higher than potential
 458 FAPAR under potential natural vegetation in a given month. However this is often a temporal effect, as
 459 when PNV FAPAR is aggregated over the whole year, most places modified by human management show
 460 actual FAPAR is lower than potential. In Western Australian cropping zones for example, crop fields
 461 have higher FAPAR during the winter growing season, but since the fields are bare for most of the year,
 462 aggregated annual PNV FAPAR is higher overall. Whilst this pattern may hold for rain-fed agriculture, in
 463 intensively irrigated areas the FAPAR of the managed vegetation can be much higher than of the PNV over
 464 the whole year, especially in arid and semi-arid areas (e.g. Nile Delta). This supplemental irrigation, plus

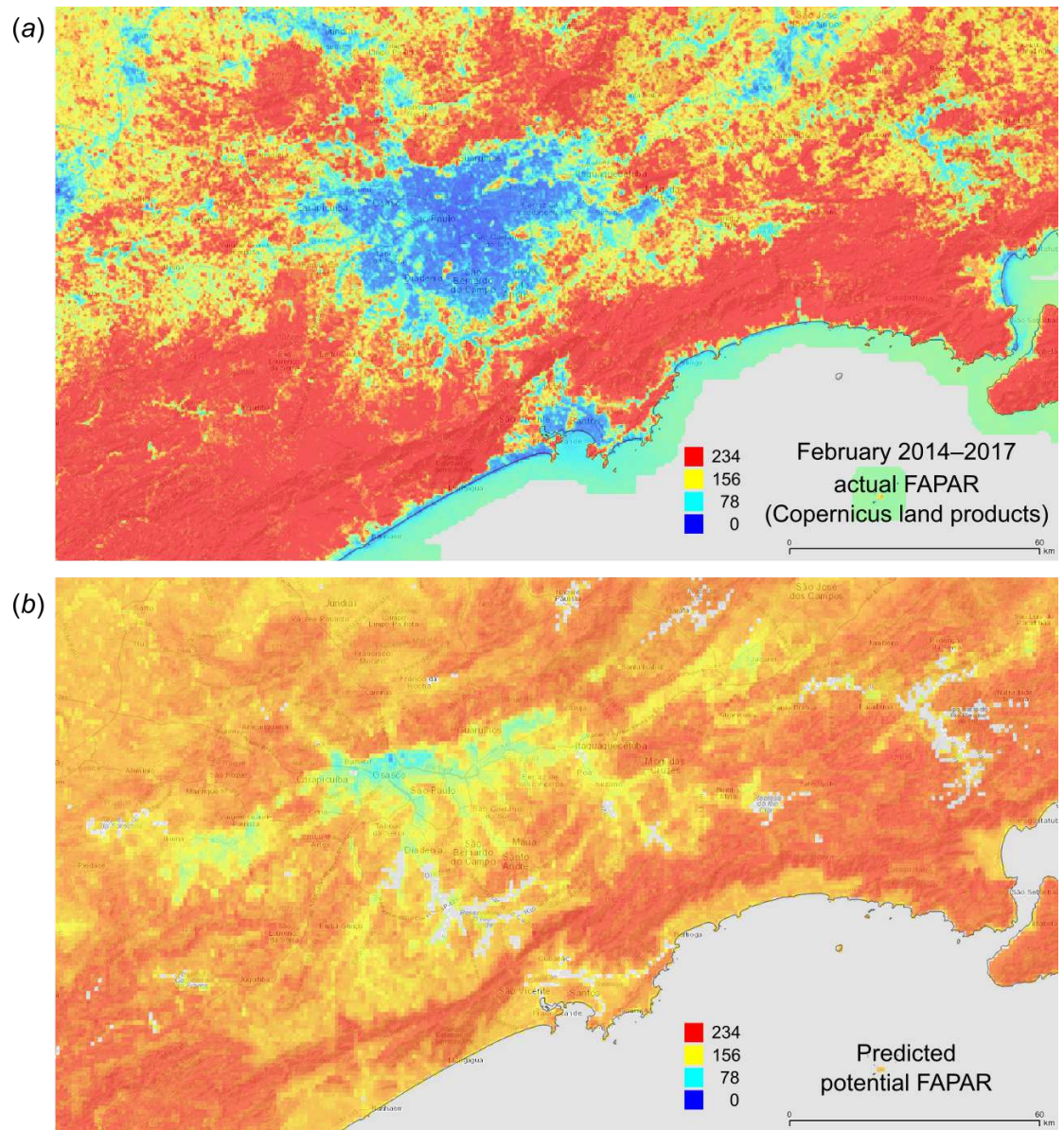


Figure 12. FAPAR values for February based on the PNV samples: (a) actual (250 m resolution) and (b) predicted (1 km resolution). A zoom in area around the city of São Paulo in Brazil.

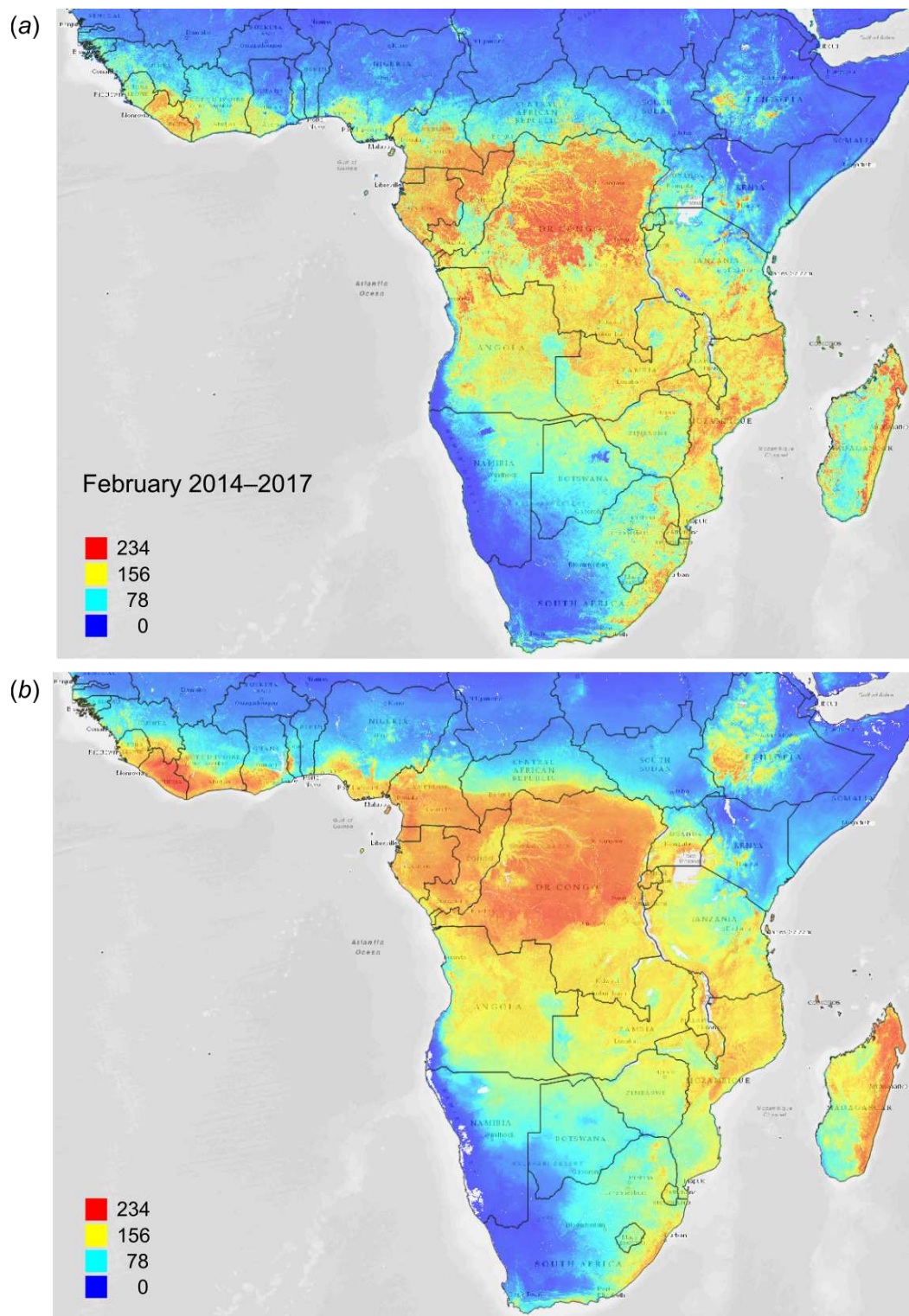


Figure 13. FAPAR values for Subsaharan Africa: (a) actual (250 m resolution) and (b) predicted (1 km resolution) potential FAPAR values for February. Background map data: Google, DigitalGlobe.

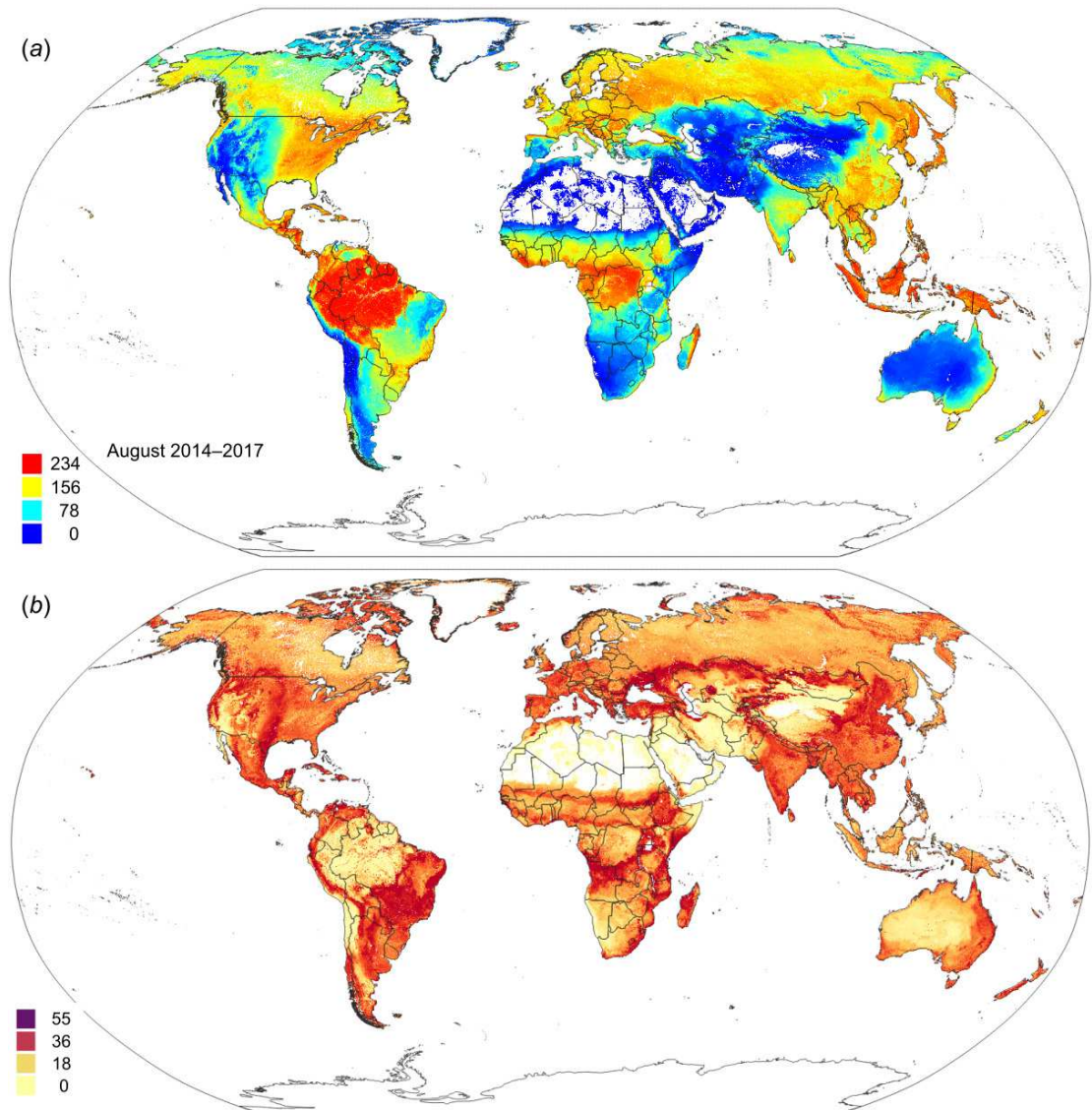


Figure 14. Predicted global FAPAR values for August (a) and standard deviation of the prediction error for the map above (b). To convert to percent divide by 253.

465 the fact that total annual precipitation is the most important covariate, indicates that water availability/use
466 efficiency is likely the main driver of FAPAR beyond natural conditions.

467 Maps of the standard deviation (s.d.) of the prediction error (Fig. 14) as derived in the ranger package
468 by using the `quantreg` setting (Meinshausen, 2006) provide useful information about model quality
469 i.e. where collection of additional points would maximize model improvement and which additional
470 covariates could be considered. For example, the highest prediction errors for FAPAR for the month of
471 August occurred in the transition areas between tropical forest and savanna areas, and in various biome
472 transition zones in Asia.

473 DISCUSSION

474 Accuracy and reliability of produced PNV maps

475 Our results of modeling potential spatial distribution of global biomes, potential FAPAR and European
476 forest tree taxa, show that relatively accurate maps of PNV can be produced using existing data and
477 publicly available environmental grids. In the case of the biomes and forest tree taxa case studies,
478 random forest consistently outperforms neural networks, gradient boosting and similar MLAs. This is
479 consistent with some other vegetation mapping studies (Li et al., 2016). However, random forest and
480 Cubist models perform equally well in the case of FAPAR. Accuracy assessment results of our work
481 indicate improvement in product accuracy in terms of greater spatial detail and smaller classification error
482 than found in the mapping products of Levvasseur et al. (2012) and Tian et al. (2016).

483 Precipitation, temperature maps and bioclimatic images are consistently the most important covariates
484 in all three case studies. Currently available lithology/parent material maps are not indicated as signifi-
485 cantly important covariates in any of the case studies. This may be because the existing lithologic map
486 (Hartmann and Moosdorf, 2012) is not detailed enough, and/or because the differences in lithology/parent
487 material are more important at finer resolutions/scales than those mapped here. Landform and lithology/-
488 parent material covariates may be important at local scales but, globally, vegetation distribution seems to
489 be dominated by climate. This is not surprising since nutrient availability is also partially controlled by
490 climate and partially by the vegetation itself. Upon visualization of the mapping products however, it was
491 noticed that the influence of topography is visible, especially in the maps of European forest tree taxa,
492 suggesting that DEM derivatives are still important for mapping PNV.

493 We have also not considered any soil layers as inputs to modeling as these are also often predicted
494 from similar climatic and remote sensing layers already used in our case studies as covariates. Moreover,
495 most of the predictive soil mapping projects use RS images reflecting human induced changes, which we
496 have tried to avoid as these are more relevant for mapping actual vegetation. For mapping of the Potential
497 Managed Vegetation, however, it would be probably more important to include also soil property / soil
498 type maps into the modeling framework.

499 Further improvements in prediction accuracy of global biome may be limited due to:

- 500 1. BIOME reconstructions representing the vegetation of an area around a given site rather than at the
501 exact point location, since the source of the pollen is on the order of 10–30 km around the site.
- 502 2. The ambiguity of reconstructions for about 10 % of the sites, so that maximum accuracy of any
503 prediction technique may not exceed 90 % without additional observation data.
- 504 3. The fact that the BIOME reconstruction accuracy is known to be lower at ecotonal boundaries and
505 in mountainous areas because of pollen transport issues, particularly the long-distance transport of
506 tree pollen.
- 507 4. The BIOME data set is compiled from many regional reconstructions and all harmonization was
508 done a posteriori, which may have introduced additional noise into the data.

509 So far, we did not explore opportunities for combining multiple MLA models based on validation
510 data i.e. for doing ensemble predictions, model averages or model stacks. Stacking models can improve
511 upon individual best techniques, achieving improvements of up to >30 %, with the additional costs
512 including higher computation loads (Michailidis, 2017). In our case, the extensive computational load
513 from derivation of models and product predictions had already obtained improved accuracies, making
514 increasing computing loads further a matter of diminishing returns.

515 Our list of MLA models could also be extended. For example, we did not consider the use of
516 Support Vector Machines (Li et al., 2016), or the Extreme Learning Machine algorithm (Deo and Şahin,
517 2015). Both have proven to be suitable for mapping vegetation distribution and quantitative properties of
518 vegetation. Not all MLA methods are, however, suitable for large regression matrices, as the computing
519 time can be excessive and hence parallelization options are crucial.

520 Our models of PNV FAPAR are based on simulated point data and the accuracy of how well models rep-
521 resent natural vegetation areas is dependent on the representativeness of the <http://protectedplanet.net>
522 and <http://intactforests.org> data. Also, many of the world's biomes such as the Mediter-
523 ranean region and similar, have sustained high levels of human impact in the past and are perhaps
524 under-represented in the <http://protectedplanet.net> data set. Nevertheless, our cross-validation
525 results (Leave-Location-Out method) indicate a good match between training and validation points.

526 It would be useful to further explore what the performance of the models we used would be if we
527 removed whole continents in the cross-validation process, or at least larger countries such as USA, China,
528 Brazil, Australia, India and/or the South African Republic. For biomes, spatial Cross Validation showed a
529 significant drop in accuracy; removing some larger countries from model training will likely also make
530 difference. We did not explore effects of spatial proximity on mapping forest species and FAPAR as these
531 are very dense point data sets. In addition, FAPAR training points were generated using simple random
532 sampling, so spatial clustering should be non-existent.

533 Fourcade et al. (2018) recently demonstrated that randomly chosen classical paintings can also be
534 added to predictive modeling, and sometimes such models might be even better evaluated than models
535 computed using real environmental variables. MLAs have even higher tendency to over-fit data and
536 often perform very poor in extrapolation areas. These two remain the biggest drawbacks of using MLAs
537 for species distribution modeling. It appears that the key to avoiding over-fitting or using non-realistic
538 mapping accuracy measures, based on Fourcade et al. (2018), is in putting more effort in cross-validation
539 (i.e. making it more robust and more reliable) and in ensuring that most important predictors and partial
540 correlations can also be explained.

541 **Possible uses of the produced PNV maps**

542 Newbold et al. (2016) argued that many terrestrial biomes today have transgressed safe limits for bio-
543 diversity, with grasslands being most affected, and tundra and boreal forests least affected. *“Slowing
544 or reversing the global loss of local biodiversity will require preserving the remaining areas of natural
545 (primary) vegetation and, so far as possible, restoring human-used lands to natural.”* (Newbold et al.,
546 2016) Roughly half of the difference of around 466 billion tonnes of carbon can be attributed to the
547 clearing of forests and woodlands, mostly for agricultural purposes (Erb et al., 2017). The other half of

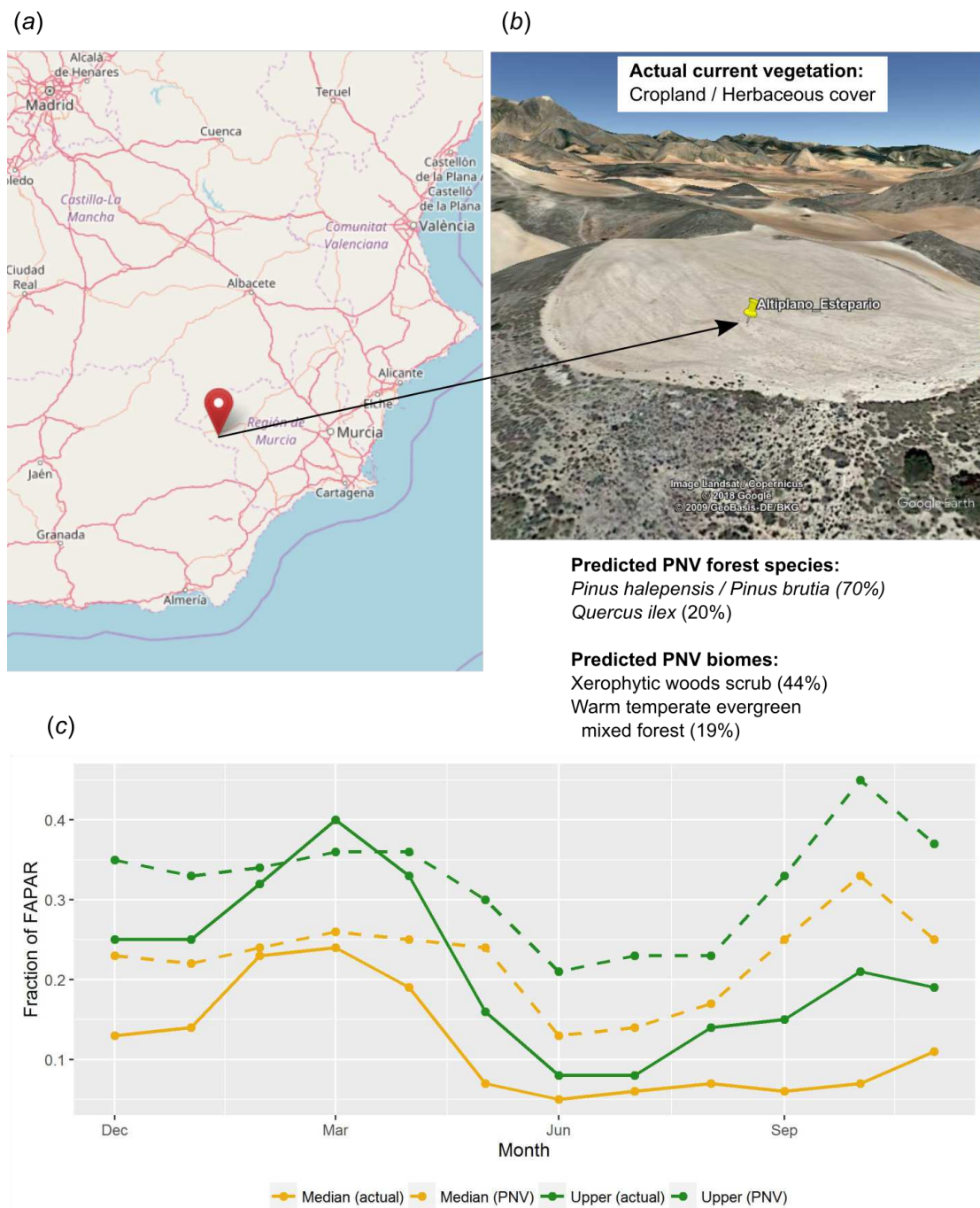


Figure 15. Example of comparison between the actual land cover and actual FAPAR curves and our predicted potential natural vegetation (PNV) and predicted PNV FAPAR curves. According to our results, this location (a–b) in southern Spain (latitude=37.938478, longitude=-2.176692) currently utilizes 51 % of the predicted FAPAR capability under PNV, indicating a substantive short fall in on-site photosynthetically active biomass (c). Background map (a) source: OpenStreetMap; landscape view (b) map data: Google, DigitalGlobe.

548 biomass carbon stock losses is derived from the management effects within a land cover class (Erb et al.,
549 2017). The expansion of agriculture will probably continue in the future, leading to decreased biodiversity

550 and soil degradation (Mauser et al., 2015; Molotoks et al., 2017). On the other hand, Griscom et al. (2017)
 551 identify reforestation (e.g. biomass restoration) as the largest natural pathway to hold global warming
 552 below 2 °C. In that context, accurate maps of PNV could become increasingly useful for assessing the
 553 level of land degradation/biomass shortfall relative to the potential of a site. Such information can also
 554 inform selection of optimal steps towards restoring biomass stocks in managed vegetation in ways that
 555 better reflect the PNV FAPAR in those areas.

556 Other uses of PNV maps include assessing the land potential i.e. land use efficiency given the difference
 557 between actual and potential vegetation. Consider for example a location in southern Spain called
 558 “*Altiplano Estepario*”, which has been identified by the Commonland company (<http://commonland.com>)
 559 and partners as a landscape restoration site. Fig. 15 shows results of a spatial query for this location
 560 and values of our PNV and PNV FAPAR predictions, in comparison to the actual land cover and actual
 561 FAPAR images. The figure shows that the actual FAPAR is as good as PNV FAPAR in February and
 562 March but that differences are large in the summer months. Overall, the median and upper FAPAR for
 563 this specific location are only 51 % of the PNV FAPAR, so we can say that this site is currently operating
 564 at 51 % of the predicted FAPAR capability under PNV. This comparison should also consider that our
 565 estimates of FAPAR come with an RMSE of ± 0.085 . Furthermore, as landscape restoration efforts have
 566 recently begun on this site — this work suggests that it ought to be possible to: (a) identify priority areas
 567 of PNV FAPAR shortfall, (b) use this information to inform in part the type of restoration strategies
 568 used, and (c) monitor the progress of restoration efforts in monthly time steps over several decades. Such
 569 practical measurement, monitoring and verification efforts are required to mobilize further investment in
 570 this emerging sector.

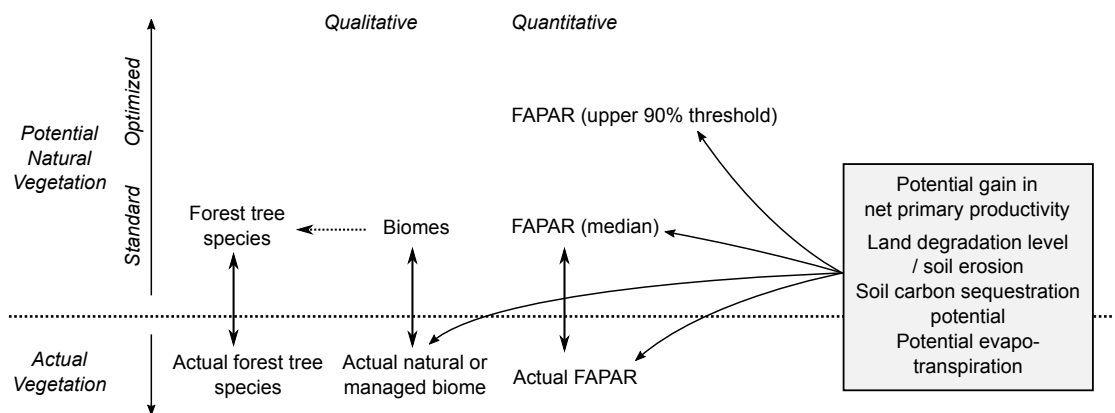


Figure 16. Some possible uses of maps of Potential Natural Vegetation.

571 Our PNV maps could also be used to estimate soil carbon sequestration and/or evapotranspiration
 572 potential, and gains in net primary productivity assuming return of natural vegetation (Fig. 16). Further
 573 more, by combining various estimates of potential natural and managed vegetation, one could design the
 574 optimal use of land both regionally and globally. Herrick et al. (2013), for example, provide a theoretical
 575 framework for estimating land potential productivity which could theoretically connect all land owners in
 576 the world to share local and regional knowledge.

577 Maps of PNV for European tree species could also be used as a supplement to the distribution and

578 ecology of tree species produced by [San-Miguel-Ayanz et al. \(2016\)](#) and [Brus et al. \(2012\)](#). Species such
579 as *Carpinus orientalis*, *Cupressus sempervirens*, *Prunus mahaleb*, *Sorbus domestica* are all predicted with
580 $TPR < 0.5$ indicating critically poor accuracy. Possible reasons for such low accuracy are problems with
581 representation of training points and somewhat too broad ecological conditions, especially if a species
582 follows some other more dominant tree species that have wide ecological niche. These maps should
583 probably not be used for spatial planning.

584 PNV for European tree species analysis could be made even more quantitative so that even predictions
585 of dendrometric properties of tree species could be produced using similar frameworks. Also, similar PNV
586 mapping algorithms could be used to map the potential canopy height based on the previously estimated
587 map of the global canopy height ([Simard et al., 2011](#)).

588 **Technical limitations and further challenges**

589 Running Machine Learning Algorithms on larger and larger data is computationally demanding; however,
590 by using fully parallelized implementation of random forest in the ranger package, we were able to
591 produce spatial predictions within days. Model fitting and prediction using EU Forest and GBIF data (1.5
592 million training points) was, however, very memory and time consuming and is not recommended for
593 systems with < 126 GiB RAM. In our case, model fitting took several hours even with full parallelization,
594 and final models were > 10 GiB in size. Prediction of probabilities took an additional 5–6 hours with the
595 current computational set-up. In the future, scalable cloud computing could be used to overcome some of
596 these computational limits. Machine learning will in any case continue to play a central role in analyzing
597 large remote sensing data stacks and extracting useful spatial patterns ([Lary et al., 2016](#)).

598 With enough computing capacity, one could theoretically use all 160 million records of distribution
599 of plant species currently available via GBIF ([Meyer et al., 2016](#)) and from other national inventories
600 to map global distribution of each forest tree species. In Europe the list is very short; globally this list
601 could be quite long (e.g. 60,000 species). The primary problems of using GBIF for PNV mapping will
602 remain however, as these are primarily due to high clustering of points and under-representation of often
603 inaccessible areas with very high biodiversity ([Yesson et al., 2007](#); [Meyer et al., 2016](#)). GBIF records have
604 been shown in the past to give biased results ([Escribano et al., 2016](#)), so that spatial prediction methods
605 that account for high spatial clustering, i.e. bias in training point representation in both space and time;
606 would need to be developed further to minimize such effects.

607 **CONCLUSIONS**

608 Although PNV is a hypothetical concept, ground-truth observations can be used to cross-validate PNV
609 models and produce an objective estimate of accuracy. As the prediction accuracy becomes more
610 significant, the reliability of the PNV maps increases. Our analyses show that the highest accuracy for
611 predicting 20 biome classes is about 68 % (33 % with spatial Cross Validation) with the most important
612 predictors being total annual precipitation, monthly temperatures and bioclimatic layers. Predictions of 73
613 forest tree species had a mapping accuracy of 25 % and with average TPR of 0.69, with the most important
614 predictors being mean annual and monthly temperatures, elevation and monthly cloud fraction. Regression
615 models for FAPAR (monthly images) were most accurate with R-square of 90 % (Leave-Location-Out CV)

616 and with the most important predictors being total annual precipitation, MODIS cloud fraction images,
617 CHELSA bioclimatic layers and month of the year, respectively. Machine learning can be successfully
618 used to model vegetation distribution, and is especially applicable when the training data sets consist of a
619 large number of observations and a large number of covariates. Extending the coverage of observations of
620 natural and managed vegetation, including through making new ground observations, will allow regular
621 improvements of such PNV maps.

622 ACKNOWLEDGMENTS

623 This research is a contribution to the AXA Chair Programme in Biosphere and Climate Impacts and
624 the Imperial College initiative on Grand Challenges in Ecosystems and the Environment (ICP). Authors
625 are grateful to Karger et al. (2017) for maintaining the CHELSA Climate images, US agencies NASA
626 and USGS for distributing high resolution images of Earth's atmosphere and the European Copernicus
627 Land program. We are grateful to Mauri et al. (2017) for sharing the EU-Forest — a high-resolution
628 tree occurrence dataset for Europe. We are also grateful to the Open Source software developers of the
629 packages ranger, xgboost, caret, raster, GDAL, SAGA GIS and similar, and without which this work
630 would have not be possible.

631 REFERENCES

- 632 Biau, G. and Scornet, E. (2016). A random forest guided tour. *TEST*, 25(2):197–227.
- 633 Bigelow, N. H., Brubaker, L. B., Edwards, M. E., Harrison, S. P., Prentice, I. C., Anderson, P. M., Andreev,
634 A. A., Bartlein, P. J., Christensen, T. R., Cramer, W., Kaplan, J. O., Lozhkin, A. V., Matveyeva, N. V.,
635 Murray, D. F., McGuire, A. D., Razzhivin, V. Y., Ritchie, J. C., Smith, B., Walker, D. A., Gajewski, K.,
636 Wolf, V., Holmqvist Björn, H., Igarashi, Y., Kremenetskii, K., Paus, A., Pisaric, M. F. J., and Volkova,
637 V. S. (2003). Climate change and Arctic ecosystems: 1. Vegetation changes north of 55 N between the
638 last glacial maximum, mid-Holocene, and present. *Journal of Geophysical Research: Atmospheres*,
639 108(D19).
- 640 Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., and Jones, Z. M.
641 (2016). mlr: Machine Learning in R. *Journal of Machine Learning Research*, 17(170):1–5.
- 642 Bohn, U., Zazanashvili, N., and Nakhutsrishvili, G. (2007). The map of the natural vegetation of europe
643 and its application in the caucasus ecoregion. *Bulletin of the Georgian National Academy of Sciences*,
644 175:112–121.
- 645 Borda, M. (2011). *Fundamentals in Information Theory and Coding*. Springer Berlin Heidelberg.
- 646 Breiman, L. (2001). Random forests. *Machine learning*, 45(1):5–32.
- 647 Brenning, A. (2012). Spatial cross-validation and bootstrap for the assessment of prediction rules in
648 remote sensing: The R package sperrorest. In *2012 IEEE International Geoscience and Remote Sensing
649 Symposium*, pages 5372–5375.
- 650 Brus, D., Hengeveld, G., Walvoort, D., Goedhart, P., Heidema, A., Nabuurs, G., and Gunia, K. (2012).
651 Statistical mapping of tree species over europe. *European Journal of Forest Research*, 131(1):145–157.

- 652 Carnahan, J. (1989). *Australia natural vegetation: Australia's vegetation in the 1780's*. Australian
653 Surveying and Land Information Group, Dept. of Administrative Services, Queensland, Australia.
- 654 Chen, T. and Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In *Proceedings of the*
655 *22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*,
656 pages 785–794, New York, NY, USA. ACM.
- 657 Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V., and
658 Böhner, J. (2015). System for automated geoscientific analyses (saga) v. 2.1.4. *Geoscientific Model*
659 *Development*, 8(7):1991–2007.
- 660 Cutler, D. R., Edwards, T. C., Beard, K. H., Cutler, A., Hess, K. T., Gibson, J., and Lawler, J. J. (2007).
661 Random forests for classification in ecology. *Ecology*, 88(11):2783–2792.
- 662 Deo, R. C. and Şahin, M. (2015). Application of the extreme learning machine algorithm for the prediction
663 of monthly effective drought index in eastern australia. *Atmospheric Research*, 153:512–525.
- 664 Elith, J. and Leathwick, J. R. (2009). Species distribution models: ecological explanation and prediction
665 across space and time. *Annual Review of Ecology, Evolution, and Systematics*, 40(1):677.
- 666 Erb, K.-H., Kastner, T., Plutzer, C., Bais, A. L. S., Carvalhais, N., Fetzel, T., Gingrich, S., Haberl, H.,
667 Lauk, C., Niedertscheider, M., Pongratz, J., Thurner, M., and Luysaert, S. (2017). Unexpectedly large
668 impact of forest management and grazing on global vegetation biomass. *Nature*, 553:73–.
- 669 Escribano, N., Ariño, A. H., and Galicia, D. (2016). Biodiversity data obsolescence and land uses changes.
670 *PeerJ*, 4:e2743.
- 671 Fan, Y., Li, H., and Miguez-Macho, G. (2013). Global patterns of groundwater table depth. *Science*,
672 339(6122):940–943.
- 673 Fluet-Chouinard, E., Lehner, B., Rebelo, L.-M., Papa, F., and Hamilton, S. K. (2015). Development of a
674 global inundation map at high spatial resolution from topographic downscaling of coarse-scale remote
675 sensing data. *Remote Sensing of Environment*, 158:348–361.
- 676 Fourcade, Y., Besnard, A. G., and Secondi, J. (2018). Paintings predict the distribution of species, or
677 the challenge of selecting environmental predictors and evaluation statistics. *Global Ecology and*
678 *Biogeography*, 27(2):245–256.
- 679 Friedman, J. H. (2002). Stochastic gradient boosting. *Computational Statistics & Data Analysis*,
680 38(4):367–378.
- 681 Griscorn, B. W., Adams, J., Ellis, P. W., Houghton, R. A., Lomax, G., Miteva, D. A., Schlesinger, W. H.,
682 Shoch, D., Siikamäki, J. V., Smith, P., Woodbury, P., Zganjar, C., Blackman, A., Campari, J., Conant,
683 R. T., Delgado, C., Elias, P., Gopalakrishna, T., Hamsik, M. R., Herrero, M., Kiesecker, J., Landis, E.,
684 Laestadius, L., Leavitt, S. M., Minnemeyer, S., Polasky, S., Potapov, P., Putz, F. E., Sanderman, J.,
685 Silvius, M., Wollenberg, E., and Fargione, J. (2017). Natural climate solutions. *Proceedings of the*
686 *National Academy of Sciences*, 114(44):11645–11650.
- 687 Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A., Thau, D.,
688 Stehman, S. V., Goetz, S. J., Loveland, T. R., Kommareddy, A., Egorov, A., Chini, L., Justice, C. O.,
689 and Townshend, J. R. G. (2013). High-resolution global maps of 21st-century forest cover change.
690 *Science*, 342(6160):850–853.
- 691 Harrison, S., Yu, G., Takahara, H., and Prentice, I. (2001). Plant diversity and palaeovegetation in East

- 692 Asia. *Nature*, 413:129–130.
- 693 Harrison, S. P. and Bartlein, P. (2012). Chapter 14 — records from the past, lessons for the future: What
694 the palaeorecord implies about mechanisms of global change. In Henderson-Sellers, A. and McGuffie,
695 K., editors, *The Future of the World's Climate (Second Edition)*, pages 403 – 436. Elsevier, Boston,
696 second edition edition.
- 697 Hartmann, J. and Moosdorf, N. (2012). The new global lithological map database GLiM: A representation
698 of rock properties at the Earth surface. *Geochemistry, Geophysics, Geosystems*, 13(12):n/a–n/a.
- 699 Hengl, T., Nussbaum, M., Wright, M. N., and Heuvelink, G. B. (2018). Random forest as a generic frame-
700 work for predictive modeling of spatial and spatio-temporal variables. *PeerJ Preprints*, 6:e26693v1.
- 701 Herrick, J. E., Urama, K. C., Karl, J. W., Boos, J., Johnson, M.-V. V., Shepherd, K. D., Hempel, J.,
702 Bestelmeyer, B. T., Davies, J., Larson Guerra, J., Kosnik, C., Kimiti, D. W., Losinyen Ekai, A.,
703 Muller, K., Norfleet, L., Ozor, N., Reinsch, T., Sarukhan, J., and West, L. T. (2013). The global
704 Land-Potential Knowledge System (LandPKS): Supporting evidence-based, site-specific land use and
705 management through cloud computing, mobile applications, and crowdsourcing. *Journal of Soil and*
706 *Water Conservation*, 68(1):5A–12A.
- 707 Hijmans, R. J. and Elith, J. (2018). *Species distribution modeling with R*. Environmental Science and
708 Policy, University of California.
- 709 Karger, D. N., Conrad, O., Böhrner, J., Kawohl, T., Kreft, H., Soria-Auza, R. W., Zimmermann, N. E.,
710 Linder, H. P., and Kessler, M. (2017). Climatologies at high resolution for the earth's land surface areas.
711 *Scientific data*, 4:170122.
- 712 Kuhn, M. (2008). Building Predictive Models in R Using the caret Package. *Journal of Statistical*
713 *Software*, 28(1):1–26.
- 714 Kuhn, M. and Johnson, K. (2013). *Applied predictive modeling*. Springer.
- 715 Kuhn, M., Weston, S., Keefer, C., Coulter, N., and Quinlan, R. (2017). *Cubist: rule-and instance-based*
716 *regression modeling*. R package version 0.2.2.
- 717 Lary, D. J., Alavi, A. H., Gandomi, A. H., and Walker, A. L. (2016). Machine learning in geosciences and
718 remote sensing. *Geoscience Frontiers*, 7(1):3 – 10. Special Issue: Progress of Machine Learning in
719 Geosciences.
- 720 Leong, M. and Roderick, G. K. (2015). Remote sensing captures varying temporal patterns of vegetation
721 between human-altered and natural landscapes. *PeerJ*, 3:e1141.
- 722 Levavasseur, G., Vrac, M., Roche, D., and Paillard, D. (2012). Statistical modelling of a new global
723 potential vegetation distribution. *Environmental Research Letters*, 7(4):044019.
- 724 Li, X., Chen, W., Cheng, X., and Wang, L. (2016). A comparison of machine learning algorithms for
725 mapping of complex surface-mined and agricultural landscapes using ziyuan-3 stereo satellite imagery.
726 *Remote sensing*, 8(6):514.
- 727 Marchant, R., Cleef, A., Harrison, S. P., Hooghiemstra, H., Markgraf, V., van Boxel, J., Ager, T., Almeida,
728 L., Anderson, R., Baied, C., Behling, H., Berrio, J. C., Burbridge, R., Björck, S., Byrne, R., Bush, M.,
729 Duivenvoorden, J., Flenley, J., De Oliveira, P., van Geel, B., Graf, K., Gosling, W. D., Harbele, S.,
730 van der Hammen, T., Hansen, B., Horn, S., Kuhry, P., Ledru, M.-P., Mayle, F., Leyden, B., Lozano-
731 Garcia, S., Melief, A. M., Moreno, P., Moar, N. T., Prieto, A., van Reenen, G., Salgado-Labouriau, M.,

- 732 Schäbitz, F., Schreve-Brinkman, E. J., and Wille, M. (2009). Pollen-based biome reconstructions for
733 latin america at 0, 6000 and 18 000 radiocarbon years ago. *Climate of the Past*, 5:725–767.
- 734 Marinova, E., Harrison, S. P., Bragg, F., Connor, S., Laet, V., Leroy, S. A., Mudie, P., Atanassova,
735 J., Bozilova, E., Caner, H., Cordova, C., Djamali, M., Filipova-Marinova, M., Gerasimenko, N.,
736 Jahns, S., Kouli, K., Kotthoff, U., Kvavadze, E., Lazarova, M., Novenko, E., Ramezani, E., Röpke,
737 A., Shumilovskikh, L., Tantâu, I., and Tonkov, S. (2018). Pollen-derived biomes in the Eastern
738 Mediterranean–Black Sea–Caspian–Corridor. *Journal of Biogeography*, 45(2):484–499.
- 739 Mauri, A., Strona, G., and San-Miguel-Ayanz, J. (2017). EU-Forest, a high-resolution tree occurrence
740 dataset for Europe. *Scientific data*, 4:160123.
- 741 Mauser, W., Klepper, G., Zabel, F., Delzeit, R., Hank, T., Putzenlechner, B., and Calzadilla, A. (2015).
742 Global biomass production potentials exceed expected future demand without the need for cropland
743 expansion. *Nature communications*, 6:8946.
- 744 Meinshausen, N. (2006). Quantile regression forests. *Journal of Machine Learning Research*, 7(Jun):983–
745 999.
- 746 Meyer, C., Weigelt, P., and Kreft, H. (2016). Multidimensional biases, gaps and uncertainties in global
747 plant occurrence information. *Ecology Letters*, 19(8):992–1006.
- 748 Michailidis, M. (2017). *Investigating machine learning methods in recommender systems*. PhD thesis,
749 UCL (University College London).
- 750 Mitchell, T. and GDAL Developers (2014). *Geospatial Power Tools: GDAL Raster & Vector Commands*.
751 Locate Press.
- 752 Molotoks, A., Kuhnert, M., Dawson, T. P., and Smith, P. (2017). Global Hotspots of Conflict Risk between
753 Food Security and Biodiversity Conservation. *Land*, 6(4).
- 754 Newbold, T., Hudson, L. N., Arnell, A. P., Contu, S., De Palma, A., Ferrier, S., Hill, S. L. L., Hoskins,
755 A. J., Lysenko, I., Phillips, H. R. P., Burton, V. J., Chng, C. W. T., Emerson, S., Gao, D., Pask-Hale, G.,
756 Hutton, J., Jung, M., Sanchez-Ortiz, K., Simmons, B. I., Whitmee, S., Zhang, H., Scharlemann, J. P. W.,
757 and Purvis, A. (2016). Has land use pushed terrestrial biodiversity beyond the planetary boundary? a
758 global assessment. *Science*, 353(6296):288–291.
- 759 Omernik, J. M. (1987). Ecoregions of the conterminous united states. *Annals of the Association of*
760 *American geographers*, 77(1):118–125.
- 761 Østbye Hemsing, L. and Bryn, A. (2012). Three methods for modelling potential natural vegetation (pnv)
762 compared: A methodological case study from south-central norway. *Norsk Geografisk Tidsskrift —*
763 *Norwegian Journal of Geography*, 66(1):11–29.
- 764 Pekel, J.-F., Cottam, A., Gorelick, N., and Belward, A. S. (2016). High-resolution mapping of global
765 surface water and its long-term changes. *Nature*, 540:418–.
- 766 Pickett, E. J., Harrison, S. P., Hope, G., Harle, K., Dodson, J. R., Peter Kershaw, A., Colin Prentice,
767 I., Backhouse, J., Colhoun, E. A., D’Costa, D., Flenley, J., Grindrod, J., Haberle, S., Hassell, C.,
768 Kenyon, C., Macphail, M., Martin, H., Martin, A. H., McKenzie, M., Newsome, J. C., Penny, D.,
769 Powell, J., Ian Raine, J., Southern, W., Stevenson, J., Sutra, J.-P., Thomas, I., Kaars, S., and Ward,
770 J. (2004). Pollen-based reconstructions of biome distributions for Australia, Southeast Asia and the
771 Pacific (SEAPAC region) at 0, 6000 and 18,000 14C yr BP. *Journal of Biogeography*, 31(9):1381–1444.

- 772 Potapov, P., Laestadius, L., and Minnemeyer, S. (2011). *Global Map of Potential Forest Cover*. World
773 Resources Institute.
- 774 Potapov, P., Yaroshenko, A., Turubanova, S., Dubinin, M., Laestadius, L., Thies, C., Aksenov, D., Egorov,
775 A., Yesipova, Y., Glushkov, I., Karpachevskiy, M., Kostikova, A., Manisha, A., Tsybikova, E., and
776 Zhuravleva, I. (2008). Mapping the world's intact forest landscapes by remote sensing. *Ecology and
777 Society*, 13(2).
- 778 Prentice, I. C. and Jolly, D. (2000). Mid-holocene and glacial-maximum vegetation geography of the
779 northern continents and africa. *Journal of Biogeography*, 27(3):507–519.
- 780 Ridgeway, G. (2017). *gbm: generalized boosted regression models*. R package version 1.6-3.1.
- 781 Ripley, B. and Venables, W. (2017). *nnet: Feed-Forward Neural Networks and Multinomial Log-Linear
782 Models*. R package version 7.3-12.
- 783 San-Miguel-Ayanz, J., de Rigo, D., Caudullo, G., Houston Durrant, T., and Mauri, A. (2016). *European
784 Atlas of forest tree species*. European Commission, Joint Research Centre.
- 785 Shannon, C. E. (1949). Communication in the presence of noise. *Proceedings of the IRE*, 37(1):10–21.
- 786 Simard, M., Pinto, N., Fisher, J. B., and Baccini, A. (2011). Mapping forest canopy height globally with
787 spaceborne lidar. *Journal of Geophysical Research: Biogeosciences*, 116(G4):NA.
- 788 Sing, T., Sander, O., Beerenwinkel, N., and Lengauer, T. (2005). ROCR: visualizing classifier performance
789 in R. *Bioinformatics*, 21(20):3940–3941.
- 790 Sing, T., Sander, O., Beerenwinkel, N., and Lengauer, T. (2016). *ROCR: Visualizing the Performance of
791 Scoring Classifiers*. R package version 1.0.7.
- 792 Tian, H., Lu, C., Ciais, P., Michalak, A. M., Canadell, J. G., Saikawa, E., Huntzinger, D. N., Gurney, K. R.,
793 Sitch, S., Zhang, B., Yang, J., Bousquet, P., Bruhwiler, L., Chen, G., Dlugokencky, E., Friedlingstein,
794 P., Melillo, J., Pan, S., Poulter, B., Prinn, R., Saunio, M., Schwalm, C. R., and Wofsy, S. C. (2016).
795 The terrestrial biosphere as a net source of greenhouse gases to the atmosphere. *Nature*, 531:225–.
- 796 Veloso, H. P., Oliveira-Filho, L., Vaz, A., Lima, M., Marquete, R., and Brazao, J. (1992). *Manual técnico
797 da vegetação brasileira*. IBGE, Rio de Janeiro.
- 798 Venables, W. N. and Ripley, B. D. (2002). *Modern applied statistics with S*. Springer-Verlag, New York,
799 4th edition.
- 800 Weisman, A. (2012). *The world without us*. Ebury Publishing.
- 801 Wilson, A. M. and Jetz, W. (2016). Remotely sensed high-resolution global cloud dynamics for predicting
802 ecosystem and biodiversity distributions. *PLOS Biology*, 14(3):1–20.
- 803 Wright, M. N. and Ziegler, A. (2016). ranger: A Fast Implementation of Random Forests for High
804 Dimensional Data in C++ and R. *Journal of Statistical Software*, page 18.
- 805 Yesson, C., Brewer, P. W., Sutton, T., Caithness, N., Pahwa, J. S., Burgess, M., Gray, W. A., White, R. J.,
806 Jones, A. C., Bisby, F. A., and Culham, A. (2007). How Global Is the Global Biodiversity Information
807 Facility? *PLoS ONE*, 2(11):e1124.