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

Tomislav Hengl $^{\rm Corresp.,~1}$, Markus G Walsh $^{2,\,3}$, Jonathan Sanderman 4 , Ichsani Wheeler 1 , Sandy P Harrison 5 , Iain C Prentice 6

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 pollenbased 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.

- Global Mapping of Potential Natural
- ² Vegetation: An Assessment of Machine
- 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 ¹Envirometrix Ltd., Wageningen, the Netherlands
- ⁸ ²The Earth Institute, Columbia University, USA
- ⁹ ³Selian Agricultural Research Inst., Arusha, Tanzania
- ¹⁰ ⁴Woods Hole Research Center, MA USA
- ¹¹ ⁵School of Archeology, Geography and Environmental Science, University of Reading,
- 12 **UK**
- ¹³ ⁶AXA 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 Hengl 1
- ¹⁸ Email address: tom.hengl@envirometrix.net

¹⁹ ABSTRACT

Peer

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

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

INTRODUCTION 43

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

- 57
- Levayasseur et al. (2012) and Tian et al. (2016) predict global PNV classes using environmental covariates 58

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

A common limitation of existing maps is their coarse spatial resolution (about 25 km) limiting the use of these maps for operational planning (e.g. Marchant et al. (2009); Levavasseur et al. (2012) and Tian et al. (2016)). In addition, comparisons of multiple overlapping sources of PNV maps shows that they rarely agree with one another since they do not share the same mapping criteria and, traditionally,

- emphasize regionally-specific botanical groupings rather than functional classifications. Limitations of
- 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.
- Here we provide an update of comparable global PNV maps produced by Potapov et al. (2011); 68 Levayasseur et al. (2012); Tian et al. (2016) and Erb et al. (2017). We explore the possibility of increasing 69 the mapping accuracy using up-to-date maps of climate, atmosphere dynamics, landform and lithology, 70 and state-of-the-art machine learning methods. Our final aim is to produce PNV maps that are more 71 detailed, richer in information, based on objective reproducible methods; and potentially more usable 72 for global modeling and awareness raising projects. We focus on improving the spatial detail, thematic 73 accuracy and reproducibility of maps, at the cost of increasing the total computing load. We also consider 74 automation of the prediction process so that the maps can be rapidly updated as new ground truth data is 75
- ⁷⁶ obtained. Our modeling follows three phases:
- (a) model selection: we compare possible models of interest for PNV mapping and choose the optimal
 spatial prediction framework based on the cross-validation results,
- (b) model assessment: we assess the uncertainty of predictions per vegetation class and try to determine
 objectively the limitations of the mapping products for wider uses, and
- (c) prediction: we use the best performing models to produce spatial predictions, then visually assess
 maps and if necessary repeat steps a–c.

METHODS AND MATERIALS

84 Theory

- PNV is the hypothetical vegetation cover that would be present if the vegetation were in equilibrium with 85 environmental controls, including climatic factors and disturbance, and not subject to human management. 86 When considering PNV, one needs to distinguish between potential "natural' and potential "managed" 87 vegetation, and "actual" natural and "actual" managed vegetation (Fig. 1a). Vegetation is in general 88 a dynamic feature. Also PNV changes as the climatic conditions change. For example, with the future 89 global warming and changes in our climate, PNV might be significantly different than pre-industrial 90 revolution. Therefore it is important to reference PNV to the time period of interest, so that historic PNV 91 and current or future PNV maps can be produced (Fig. 1b). 92 In addition to the differentiation between the potential and actual natural vegetation, there are also 93
- ⁹⁴ three sub-types of the PNV that need to be considered:
- 1. PNV model A: based on the autochthonous or native vegetation and living species only.

Manuscript to be reviewed



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.

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

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

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

- Regardless of the hypothetical nature of PNV, the concept (both as a classification and as a regression type problem) is still a helpful yardstick against which land cover change can be quantitatively measured and land restoration designs can be planned. Erb et al. (2017) have estimated that almost half of the standing global vegetation biomass carbon stocks has been lost, almost equally due to land cover change (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

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
- modeling of the relationship between species (or natural associations of species or communities) and a

Manuscript to be reviewed

- ¹²¹ 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
- 124 and/or biophysical conditions:

$$Pr(Y) = f(Relief, BioClimate, Lithology)$$
(1)

- where Y is the target variable, which could be vegetation types or plant species with a finite number of
- states $Y \in \{1, 2, ..., k\}$ 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
- ¹²⁸ continuous vegetation properties such as biomass or leaf area index.
- ¹²⁹ The primary assumptions we make when applying a PNV model to the training data are:
- 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.
- 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:

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

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

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

- ¹⁴⁶ 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

Manuscript to be reviewed

- 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
- (Elith and Leathwick, 2009).

153 Input data: training points

- ¹⁵⁴ We consider three ground-truth data sets for model calibration:
- 1. an expanded version of the BIOME 6000 DB data set representing site-based reconstructions from
- surface pollen samples of major vegetation types or biomes (http://dx.doi.org/10.17864/
 1947.99),
- EU Forest (Mauri et al., 2017) and GBIF (Global Biodiversity Information Facilities) occurrence
 records of the 76 main forest tree taxa in Europe (http://dx.doi.org/10.15468/dl.fhucwx),
- Long-term Fraction of Absorbed Photosynthetically Active Radiation (FAPAR) monthly images
 derived using a time-series of Copernicus Global Land products (https://land.copernicus.
 eu),
- ¹⁶³BIOME 6000 and EU Forest and GBIF occurrences are point data sets, while FAPAR consists of ¹⁶⁴remote sensing images at relatively fine spatial resolution (250 m), from which we sample a large number ¹⁶⁵of values (ca 100,000) using random sampling after masking for areas of natural vegetation.

166 BIOME 6000

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

- The BIOME 6000 data set is an amalgamation of multiple data sets. BIOME 6000 initially produced 174 maps for individual regions: Europe, Africa and the Arabian Peninsula, the Former Soviet Union and 175 Mongolia and China. Additional regions were subsequently added including Beringia, western North 176 America, Canada and the eastern United States and Japan, and the data for northern Eurasia, China 177 178 and southern Europe and Africa were also updated. These regional compilations were summarized in Prentice and Jolly (2000). Subsequent regional updates include China (Harrison et al., 2001), the 179 circum-Artic region (Bigelow et al., 2003), Australia (Pickett et al., 2004) and South America (Marchant 180 et al., 2009). Additionally, we have also combined these data with pollen-based vegetation reconstructions 181
- 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.

Manuscript to be reviewed



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

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

The number of biomes differentiated varies from region to region, and some biomes were only 190 reconstructed in specific regions where they are particularly characteristic, although they may occur, but 191 not be recognized, elsewhere. Furthermore, some biomes that can be recognized on the modern landscape 192 were never reconstructed in the BIOME 6000 data set (e.g. cushion forb tundra) — either because of 193 the sample distribution or because the characteristic plant-functional types were also spread amongst 194 other biomes. Simplified or "megabiome" classifications (e.g. Harrison and Bartlein (2012)) involve a 195 substantial loss of information. We have therefore created a new standardization of the classification 196 scheme (see further Table 1; the final scheme has 20 globally applicable and distinctive biomes) which 197 preserve the maximum number of distinct biomes that were reconstructed as present in multiple regions. 198 There are relatively few data vegetation reconstructions for tropical South America, which could lead 190 to extrapolation problems and omission of important PNV classes in Latin America, but also potentially in 200 201 tropical parts of Africa and Asia. To reduce under-representation of tropics, we have added 350 randomly simulated points based on the RADAM Brazil natural vegetation polygon map at high spatial detail 202 (Radam Vegetação SIRGAS map) (Veloso et al., 1992) obtained from ftp://geoftp.ibge.gov.br/. 203 Before generating the pseudo-observations for Brazil, we translated SIRGAS map legends to match the 204 BIOME 6000 classes. This translation is also available via the project's github repository. This gave a 205 total of 8057 unique individual locations represented in the combined data set i.e. a total of 8959 training 206 observations (Fig. 2). 207

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

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

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

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

Manuscript to be reviewed

237 FAPAR

- 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
- factor 1/255 i.e. with a maximum value of 0.94). From a total of 142 images downloaded from https:
- ²⁴¹ //land.copernicus.eu we derived minimum, median and maximum value of FAPAR per month (12)
- using the 95 % probability interval using the data.table package (http://r-datatable.com). For
- ²⁴³ regression modeling we only report results of predictions of median values of FAPAR; predictions of
- ²⁴⁴ minimum and maximum FAPAR can be obtained from the data repository.



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.

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

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

R> FAPAR ~ cm + X1m + X2m + X3 + ... + Xp

²⁵³ where X1m is the covariate with monthly values (for example precipitation, day-time and night-time

temperatures etc), X3 is the environmental covariates that do not vary through year (e.g. lithology or DEM

derivatives), and cm is the cosine of the month number:

$$c_m = \cos\left(\mu/12 \cdot 2 \cdot \pi\right) \tag{4}$$

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).

- For PNV FAPAR mapping we have masked out all water bodies including lakes and rivers, following 258 the ESA's global land cover maps for the period 2000–2015 (https://www.esa-landcover-cci.org) 259 and permanent ice/barren ground. 260 Input data: environmental covariates 261 For modeling purposes, we use a stack of 160 spatially explicit co-variate data layers that represent 262 standard ecological gradients essential for growth and survival of plants: 263 • DEM derivatives quantifying various landscape metrics and hydrological processes: slope, curva-264 ture, topographic index, topographic openness, valley depth and multi-resolution valley bottom 265 index; all derived using the SAGA GIS (Conrad et al., 2015); 266 Mean, minimum and maximum monthly temperatures derived as a mean between WorldClim v2 267 (http://worldclim.org/version2) and CHELSA climate (Karger et al., 2017). 268 • Mean monthly precipitation images derived as a weighted average between the WorldClim v2, 269 CHELSA climate and Global Precipitation Measurement Integrated Multi-satellitE Retrievals for 270 GPM (IMERG) rainfall product. 271 CHELSA Bioclimatic layers downloaded from http://chelsa-climate.org/, including: an-272 nual mean temperature, mean diurnal temperature range, isothermality (day-to-night temperature 273 oscillations relative to the summer-to-winter oscillations), temperature seasonality (standard de-274 viation of monthly temperature averages), maximum temperature of warmest month, minimum 275 temperature of coldest month, temperature annual range, mean temperature of warmest quarter, 276 mean temperature of coldest quarter, annual precipitation amount, precipitation of wettest month, 277 precipitation of driest month, precipitation of wettest quarter, precipitation of driest quarter (Karger 278 et al., 2017); 279 · European Space Agency's CCI-LC snow probability monthly averages based on MODIS snow 280 products MOD10A2 downloaded from http://maps.elie.ucl.ac.be/CCI/viewer/index. 281 php; 282 USGS Global Ecophysiography landform classification and lithological map at 250 m resolution 283 obtained from http://rmgsc.cr.usgs.gov/outgoing/ecosystems/Global/ and based on 284 Global Lithological Map (GLiM) (Hartmann and Moosdorf, 2012); 285 • MODIS Cloud fraction monthly images obtained from http://www.earthenv.org/cloud (Wil-286 son and Jetz, 2016); 287 • Global Water Table Depth in meters based on Fan et al. (2013); 288 NASA's monthly MODIS Precipitable Water Vapor images (MYDAL2_M_SKY_WV data set at http: 289 //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 (Pekel et al., 2016); 294 • Density of earthquakes based on the USGS Earthquake Archives (http://earthquake.usgs. 295
- 296 gov/earthquakes/);

Manuscript to be reviewed

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

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

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

311 Machine Learning Algorithms (MLA) examined

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

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

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

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

- xgboost package (Chen and Guestrin, 2016).
- 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

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

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)</pre>

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

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

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

355 Performance of classification algorithms

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

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

369 (Shannon, 1949; Borda, 2011):

$$SSEI_{s}(x) = -\sum_{i=1}^{b} P_{i}(x) \cdot \log_{b} P_{i}(x) = \frac{-\sum_{i=1}^{b} P_{i}(x) \cdot \log P_{i}(x)}{-b \cdot b^{-1} \cdot \log b^{-1}}$$
(5)

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

For the biomes data set, where spatial clustering of points is significant, we also use repeated spatial 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)
```

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

We also compare results of modeling potential distribution of tree species in Europe with the habitat

- 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

Performance of regression algorithms is also assessed using 5–fold cross-validation with refitting of models. For assessment of the mapping accuracy for FAPAR we use as the main performance measures

³⁸⁷ the root mean squared error (RMSE):

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

³⁸⁸ and mean error (ME):

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

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

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

13/35

- 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

Results showed that a relatively accurate model of PNV could be produced from the BIOME 6000 data set using the existing stack of covariates at 1 km spatial resolution. Results of cross-validation show the random forest (RF) model to be the best performing method and distinctively superior to all other approaches (Fig. 5). The choice of the random forest mtry parameter had little impact on overall accuracy, most likely because there was a high overlap in covariate maps so that even with smaller mtry bagging the performance was relatively similar. The best prediction accuracy from among the four methods used 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.

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

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

Manuscript to be reviewed





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	Ν
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 Levavasseur et al. (2012), except in our case the Rocky Mountains in USA and mountains
- 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
- the "cold evergreen needleleaf forest" and "cool evergreen needleleaf forest" (e.g. Scotland).
- Results of the accuracy assessment based on the spatial Cross-Validation (mlr package implementation

- (Bischl 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
- individual taxa, however, can be relatively high with TPR values of between 0.16–0.90 and an average
- value of around 0.69 (Table 2). The final maps (Fig. 9) showed a relatively good match with ground
- data, meaning that with the exception of some species of rarer occurrence (*Picea omorika*, *Cupressus*
- 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	Ν
Abies alba	2685484	-0.01	0.77	0.92	16,150
Acer campestre	3189863	-0.01	0.65	0.83	19,819
Acer platanoides	3189846	-0.02	0.68	0.82	30,801
Acer pseudoplatanus	3189870	-0.01	0.69	0.79	65,039
Aesculus hippocastanum	3189815	-0.01	0.59	0.85	8,088
Ailanthus altissima	3190653	0.04	0.69	0.92	1,576
Alnus cordata	2876607	0.05	0.73	0.95	904
Alnus glutinosa	2876213	0.00	0.71	0.77	91,292

Continued on next page

Manuscript to be reviewed

Table 2 – Continued from previous page Species nome CRIE toxon ID ME_TPP_AUC							
Species name	GBIF taxon ID	ME	IPK	AUC	N		
Alnus incana	2876388	-0.03	0.76	0.95	6,873		
Betula spp.	2875008	-0.03	0.63	0.83	7,313		
Carpinus betulus	2875818	0.00	0.75	0.89	22,765		
Carpinus orientalis	2875780	0.07	0.21	0.92	284		
Castanea sativa Caltia guatualia	5555294	0.00	0.74	0.91	13,049		
Cettis australis	2984492	-0.01	0.54	0.92	594 807		
Cornus mas	3082203	0.05	0.51	0.90	027 0 027		
Cornus sanguinea	2875070	-0.03	0.59	0.82	0,057		
Cupressus sempervirens	2684030	-0.02	0.07	0.70	40,140		
Fuonymus europaeus	3169131	-0.04	0.61	0.70	12 119		
Fagus sylvatica	2882316	0.02	0.01	0.85	89 044		
Frangula alnus	3039454	-0.02	0.75	0.86	26 873		
Fraxinus angustifolia	7325877	-0.05	0.63	0.94	1.757		
Fraxinus excelsior	3172358	0.00	0.67	0.74	91.111		
Fraxinus ornus	3172347	0.02	0.86	0.99	2,765		
Ilex aquifolium	5414222	-0.01	0.66	0.82	26,873		
Juglans regia	3054368	-0.03	0.60	0.89	3,643		
Juniperus communis	2684709	-0.03	0.71	0.86	21,189		
Juniperus oxycedrus	2684451	-0.07	0.71	0.97	1,705		
Juniperus phoenicea	2684640	-0.07	0.74	0.98	1,137		
Juniperus thurifera	2684528	-0.03	0.87	0.99	1,886		
Larix decidua	2686212	-0.01	0.71	0.89	15,581		
Olea europaea	5415040	0.00	0.90	0.99	7,080		
Ostrya carpinifolia	5332305	0.06	0.90	0.99	1,809		
Picea abies	5284884	0.02	0.76	0.86	122,713		
Picea sitchensis	5284827	0.05	0.80	0.96	13,023		
Pinus cembra	5285134	-0.01	0.77	0.96	853		
Pinus halepensis and Pinus brutia	5285604	0.03	0.86	0.99	16,951		
Pinus mugo	5285385	0.00	0.85	0.98	6,667		
Pinus nigra Binus pinaster	5284809	0.01	0.79	0.93	13,540		
Pinus pinas	5285165	0.01	0.80	0.98	17,080		
I mus pineu Pinus sylvastris	5285637	-0.04	0.85	0.99	153 028		
Populus alba	3040233	-0.02	0.78	0.85	4 522		
Populus nigra	3040227	-0.01	0.54	0.80	5 478		
Populus tremula	3040249	-0.02	0.65	0.07	44 057		
Prunus avium	3020791	-0.01	0.63	0.77	25.711		
Prunus cerasifera	3021730	0.00	0.73	0.94	3.928		
Prunus mahaleb	3022789	-0.01	0.31	0.75	517		
Prunus padus	3021037	-0.03	0.63	0.78	21,705		
Prunus spinosa	3023221	-0.01	0.69	0.81	31,783		
Quercus cerris	2880580	0.00	0.80	0.97	4,109		
Quercus ilex	2879098	0.02	0.85	0.99	22,972		
Quercus pubescens	2881283	0.01	0.86	0.98	9,096		
Quercus pyrenaica	2878826	0.00	0.88	0.99	6,253		
Quercus robur and Quercus petraea	2878688	0.01	0.69	0.76	141,938		
Quercus suber	2879411	-0.04	0.86	0.99	5,504		
Robinia pseudoacacia	5352251	0.01	0.71	0.90	13,411		
Salix alba Salix anna a	5372513	0.02	0.72	0.90	11,938		
Salix caprea	22/2922	-0.03	0.08	0.78	40,879		
Samoucus migra Sorbus aria	2000/20	0.00	0.70	0.81	44,901 5 106		
Sorbus ana Sorbus angunaria	3012060	-0.01	0.39	0.87	2,420 86.077		
Sorbus danestica	3012107	-0.01	0.70	0.70	80,977		
Sorbus torminalis	3012567	-0.03	0.62	0.07	2 558		
Taxus baccata	5284517	-0.02	0.58	0.82	8,062		
Tilia spp.	3152041	-0.02	0.50	0.82	4.393		
Ulmus spp.	2984510	-0.03	0.64	0.92	5.426		
Tilia spp.	3152041	0.00	0.58	0.85	4,522		
Ulmus spp.	2984510	-0.02	0.69	0.91	5,375		

TT 1 1 0	<i>a</i>	10	•	
Table 2 -	- Continued	t from	previous	pag

Manuscript to be reviewed



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

Peer.

The most important predictors in the random forest model for forest tree taxa were mean annual daily
 temperature, other monthly temperatures, elevation, CHELSA bioclimatic images, monthly precipitation
 and MODIS cloud fraction images. Covariates for lithology and landform classification did not feature in
 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.

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

Manuscript to be reviewed

- 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
- 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

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



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.

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



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.

Manuscript to be reviewed

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

Our results of modeling potential spatial distribution of global biomes, potential FAPAR and European 475 forest tree taxa, show that relatively accurate maps of PNV can be produced using existing data and 476 publicly available environmental grids. In the case of the biomes and forest tree taxa case studies, 477 random forest consistently outperforms neural networks, gradient boosting and similar MLA's. This is 478 consistent with some other vegetation mapping studies (Li et al., 2016). However, random forest and 479 Cubist models perform equally well in the case of FAPAR. Accuracy assessment results of our work 480 indicate improvement in product accuracy in terms of greater spatial detail and smaller classification error 481 than found in the mapping products of Levavasseur et al. (2012) and Tian et al. (2016). 482 Precipitation, temperature maps and bioclimatic images are consistently the most important covariates 483 in all three case studies. Currently available lithology/parent material maps are not indicated as signifi-484 cantly important covariates in any of the case studies. This may be because the existing lithologic map 485 (Hartmann and Moosdorf, 2012) is not detailed enough, and/or because the differences in lithology/parent 486 material are more important at finer resolutions/scales than those mapped here. Landform and lithology/-487 parent material covariates may be important at local scales but, globally, vegetation distribution seems to 488 be dominated by climate. This is not surprising since nutrient availability is also partially controlled by 489 climate and partially by the vegetation itself. Upon visualization of the mapping products however, it was 490 noticed that the influence of topography is visible, especially in the maps of European forest tree taxa, 491

⁴⁹² suggesting that DEM derivatives are still important for mapping PNV.

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

- ⁴⁹⁹ Further improvements in prediction accuracy of global biome may be limited due to:
- BIOME reconstructions representing the vegetation of an area around a given site rather than at the
 exact point location, since the source of the pollen is on the order of 10–30 km around the site.
- 2. The ambiguity of reconstructions for about 10% of the sites, so that maximum accuracy of any
 prediction technique may not exceed 90% without additional observation data.
- The fact that the BIOME reconstruction accuracy is known to be lower at ecotonal boundaries and
 in mountainous areas because of pollen transport issues, particularly the long-distance transport of
 tree pollen.
- 4. The BIOME data set is compiled from many regional reconstructions and all harmonization was
 done a posteriori, which may have introduced additional noise into the data.

Manuscript to be reviewed

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

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

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

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

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

541 Possible uses of the produced PNV maps

Newbold et al. (2016) argued that many terrestrial biomes today have transgressed safe limits for biodiversity, with grasslands being most affected, and tundra and boreal forests least affected. "Slowing

- 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

Manuscript to be reviewed



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

Manuscript to be reviewed

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

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



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

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

Manuscript to be reviewed

- ⁵⁷⁸ ecology of tree species produced by San-Miguel-Ayanz et al. (2016) and Brus et al. (2012). Species such 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.
- PNV for European tree species analysis could be made even more quantitative so that even predictions of dendrometric properties of tree species could be produced using similar frameworks. Also, similar PNV mapping algorithms could be used to map the potential canopy height based on the previously estimated map of the global canopy height (Simard et al., 2011).

588 Technical limitations and further challenges

Running Machine Learning Algorithms on larger and larger data is computationally demanding; however, 589 by using fully parallelized implementation of random forest in the ranger package, we were able to 590 produce spatial predictions within days. Model fitting and prediction using EU Forest and GBIF data (1.5 591 million training points) was, however, very memory and time consuming and is not recommended for 592 systems with < 126 GiB RAM. In our case, model fitting took several hours even with full parallelization, 503 and final models were >10 GiB in size. Prediction of probabilities took an additional 5–6 hours with the 594 current computational set-up. In the future, scalable cloud computing could be used to overcome some of 595 these computational limits. Machine learning will in any case continue to play a central role in analyzing 596 large remote sensing data stacks and extracting useful spatial patterns (Lary et al., 2016). 597

With enough computing capacity, one could theoretically use all 160 million records of distribution 598 of plant species currently available via GBIF (Meyer et al., 2016) and from other national inventories 500 to map global distribution of each forest tree species. In Europe the list is very short; globally this list 600 could be quite long (e.g. 60,000 species). The primary problems of using GBIF for PNV mapping will 601 remain however, as these are primarily due to high clustering of points and under-representation of often 602 inaccessible areas with very high biodiversity (Yesson et al., 2007; Meyer et al., 2016). GBIF records have 603 been shown in the past to give biased results (Escribano et al., 2016), so that spatial prediction methods 604 that account for high spatial clustering, i.e. bias in training point representation in both space and time; 605 would need to be developed further to minimize such effects. 606

607 CONCLUSIONS

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

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

622 ACKNOWLEDGMENTS

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

631 **REFERENCES**

- Biau, G. and Scornet, E. (2016). A random forest guided tour. TEST, 25(2):197-227.
- Bigelow, N. H., Brubaker, L. B., Edwards, M. E., Harrison, S. P., Prentice, I. C., Anderson, P. M., Andreev,
- A. A., Bartlein, P. J., Christensen, T. R., Cramer, W., Kaplan, J. O., Lozhkin, A. V., Matveyeva, N. V.,
- Murray, D. F., McGuire, A. D., Razzhivin, V. Y., Ritchie, J. C., Smith, B., Walker, D. A., Gajewski, K.,
- 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
- last glacial maximum, mid-Holocene, and present. *Journal of Geophysical Research: Atmospheres*,
- 639 108(D19).
- Bischl, B., Lang, M., Kotthoff, L., Schiffner, J., Richter, J., Studerus, E., Casalicchio, G., and Jones, Z. M.
- (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
- 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
- remote sensing: The R package sperrorest. In 2012 IEEE International Geoscience and Remote Sensing
- ⁶⁴⁹ *Symposium*, pages 5372–5375.
- 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.

Peer

Manuscript to be reviewed

- Carnahan, J. (1989). Australia natural vegetation: Australia's vegetation in the 1780's. Australian 652 Surveying and Land Information Group, Dept. of Administrative Services, Queensland, Australia. 653
- Chen, T. and Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In Proceedings of the 654
- 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, 655
- pages 785-794, New York, NY, USA. ACM. 656
- Conrad, O., Bechtel, B., Bock, M., Dietrich, H., Fischer, E., Gerlitz, L., Wehberg, J., Wichmann, V., and 657
- Böhner, J. (2015). System for automated geoscientific analyses (saga) v. 2.1.4. Geoscientific Model 658
- Development, 8(7):1991-2007. 659
- Cutler, D. R., Edwards, T. C., Beard, K. H., Cutler, A., Hess, K. T., Gibson, J., and Lawler, J. J. (2007). 660
- Random forests for classification in ecology. *Ecology*, 88(11):2783–2792. 661
- Deo, R. C. and Şahin, M. (2015). Application of the extreme learning machine algorithm for the prediction 662
- of monthly effective drought index in eastern australia. Atmospheric Research, 153:512–525. 663
- Elith, J. and Leathwick, J. R. (2009). Species distribution models: ecological explanation and prediction 664
- across space and time. Annual Review of Ecology, Evolution, and Systematics, 40(1):677. 665
- Erb, K.-H., Kastner, T., Plutzar, C., Bais, A. L. S., Carvalhais, N., Fetzel, T., Gingrich, S., Haberl, H., 666
- Lauk, C., Niedertscheider, M., Pongratz, J., Thurner, M., and Luyssaert, S. (2017). Unexpectedly large 667
- impact of forest management and grazing on global vegetation biomass. Nature, 553:73-. 668
- Escribano, N., Ariño, A. H., and Galicia, D. (2016). Biodiversity data obsolescence and land uses changes. 669
- PeerJ, 4:e2743. 670

680

- Fan, Y., Li, H., and Miguez-Macho, G. (2013). Global patterns of groundwater table depth. Science, 671 339(6122):940-943. 672
- Fluet-Chouinard, E., Lehner, B., Rebelo, L.-M., Papa, F., and Hamilton, S. K. (2015). Development of a 673
- global inundation map at high spatial resolution from topographic downscaling of coarse-scale remote 674
- sensing data. Remote Sensing of Environment, 158:348-361. 675
- Fourcade, Y., Besnard, A. G., and Secondi, J. (2018). Paintings predict the distribution of species, or 676
- the challenge of selecting environmental predictors and evaluation statistics. Global Ecology and 677
- Biogeography, 27(2):245-256. 678
- Friedman, J. H. (2002). Stochastic gradient boosting. Computational Statistics & Data Analysis, 679 38(4):367-378.
- Griscom, B. W., Adams, J., Ellis, P. W., Houghton, R. A., Lomax, G., Miteva, D. A., Schlesinger, W. H., 681
- Shoch, D., Siikamäki, J. V., Smith, P., Woodbury, P., Zganjar, C., Blackman, A., Campari, J., Conant, 682
- R. T., Delgado, C., Elias, P., Gopalakrishna, T., Hamsik, M. R., Herrero, M., Kiesecker, J., Landis, E., 683
- Laestadius, L., Leavitt, S. M., Minnemeyer, S., Polasky, S., Potapov, P., Putz, F. E., Sanderman, J., 684
- Silvius, M., Wollenberg, E., and Fargione, J. (2017). Natural climate solutions. Proceedings of the 685
- National Academy of Sciences, 114(44):11645–11650. 686
- Hansen, M. C., Potapov, P. V., Moore, R., Hancher, M., Turubanova, S. A., Tyukavina, A., Thau, D., 687
- Stehman, S. V., Goetz, S. J., Loveland, T. R., Kommareddy, A., Egorov, A., Chini, L., Justice, C. O., 688
- and Townshend, J. R. G. (2013). High-resolution global maps of 21st-century forest cover change. 689 Science, 342(6160):850-853. 690
- Harrison, S., Yu, G., Takahara, H., and Prentice, I. (2001). Plant diversity and palaeovegetation in East 691

- ⁶⁹² Asia. *Nature*, 413:129–130.
- Harrison, S. P. and Bartlein, P. (2012). Chapter 14 records from the past, lessons for the future: What
- 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,
- second edition edition.
- Hartmann, J. and Moosdorf, N. (2012). The new global lithological map database GLiM: A representation
 of rock properties at the Earth surface. *Geochemistry, Geophysics, Geosystems*, 13(12):n/a–n/a.
- Hengl, T., Nussbaum, M., Wright, M. N., and Heuvelink, G. B. (2018). Random forest as a generic frame-
- 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.,
- 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., Böhner, J., Kawohl, T., Kreft, H., Soria-Auza, R. W., Zimmermann, N. E.,
- 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.
- Kuhn, M., Weston, S., Keefer, C., Coulter, N., and Quinlan, R. (2017). *Cubist: rule-and instance-based 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
- 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
- between human-altered and natural landscapes. *PeerJ*, 3:e1141.
- Levavasseur, G., Vrac, M., Roche, D., and Paillard, D. (2012). Statistical modelling of a new global
- 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
- mapping of complex surface-mined and agricultural landscapes using ziyuan-3 stereo satellite imagery.
 Remote sensing, 8(6):514.
- ⁷²⁷ Marchant, R., Cleef, A., Harrison, S. P., Hooghiemstra, H., Markgraf, V., van Boxel, J., Ager, T., Almeida,
- L., Anderson, R., Baied, C., Behling, H., Berrio, J. C., Burbridge, R., Björck, S., Byrne, R., Bush, M.,
- Duivenvoorden, J., Flenley, J., De Oliveira, P., van Geel, B., Graf, K., Gosling, W. D., Harbele, S.,
- van der Hammen, T., Hansen, B., Horn, S., Kuhry, P., Ledru, M.-P., Mayle, F., Leyden, B., Lozano-
- Garcia, S., Melief, A. M., Moreno, P., Moar, N. T., Prieto, A., van Reenen, G., Salgado-Labouriau, M.,

Manuscript to be reviewed

- ⁷³² Schäbitz, F., Schreve-Brinkman, E. J., and Wille, M. (2009). Pollen-based biome reconstructions for
- 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,
- J., Bozilova, E., Caner, H., Cordova, C., Djamali, M., Filipova-Marinova, M., Gerasimenko, N.,
- Jahns, S., Kouli, K., Kotthoff, U., Kvavadze, E., Lazarova, M., Novenko, E., Ramezani, E., Röpke,
- 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
- dataset for Europe. *Scientific data*, 4:160123.
- ⁷⁴¹ Mauser, W., Klepper, G., Zabel, F., Delzeit, R., Hank, T., Putzenlechner, B., and Calzadilla, A. (2015).
- Global biomass production potentials exceed expected future demand without the need for cropland
- expansion. *Nature communications*, 6:8946.
- Meinshausen, N. (2006). Quantile regression forests. *Journal of Machine Learning Research*, 7(Jun):983–
 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.
- Molotoks, A., Kuhnert, M., Dawson, T. P., and Smith, P. (2017). Global Hotspots of Conflict Risk between
- Food Security and Biodiversity Conservation. *Land*, 6(4).
- Newbold, T., Hudson, L. N., Arnell, A. P., Contu, S., De Palma, A., Ferrier, S., Hill, S. L. L., Hoskins,
- 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.,
- 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 —
- Norwegian Journal of Geography, 66(1):11–29.
- 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–.
- Pickett, E. J., Harrison, S. P., Hope, G., Harle, K., Dodson, J. R., Peter Kershaw, A., Colin Prentice,
- I., Backhouse, J., Colhoun, E. A., D'Costa, D., Flenley, J., Grindrod, J., Haberle, S., Hassell, C.,
- Kenyon, C., Macphail, M., Martin, H., Martin, A. H., McKenzie, M., Newsome, J. C., Penny, D.,
- Powell, J., Ian Raine, J., Southern, W., Stevenson, J., Sutra, J.-P., Thomas, I., Kaars, S., and Ward,
- J. (2004). Pollen-based reconstructions of biome distributions for Australia, Southeast Asia and the
- Pacific (SEAPAC region) at 0, 6000 and 18,000 14C yr BP. Journal of Biogeography, 31(9):1381–1444.

Potapov, P., Laestadius, L., and Minnemeyer, S. (2011). Global Map of Potential Forest Cover. World

773 Resources Institute.

- Potapov, P., Yaroshenko, A., Turubanova, S., Dubinin, M., Laestadius, L., Thies, C., Aksenov, D., Egorov,
- A., Yesipova, Y., Glushkov, I., Karpachevskiy, M., Kostikova, A., Manisha, A., Tsybikova, E., and
- Zhuravleva, I. (2008). Mapping the world's intact forest landscapes by remote sensing. *Ecology and*
- *Society*, 13(2).
- Prentice, I. C. and Jolly, D. (2000). Mid-holocene and glacial-maximum vegetation geography of the
- northern continents and africa. *Journal of Biogeography*, 27(3):507–519.
- Ridgeway, G. (2017). gbm: generalized boosted regression models. R package version 1.6-3.1.
- Ripley, B. and Venables, W. (2017). nnet: Feed-Forward Neural Networks and Multinomial Log-Linear
- 782 *Models*. R package version 7.3-12.
- San-Miguel-Ayanz, J., de Rigo, D., Caudullo, G., Houston Durrant, T., and Mauri, A. (2016). *European 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
- reason 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.
- 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,
- P., Melillo, J., Pan, S., Poulter, B., Prinn, R., Saunois, 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–.
- 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.
- 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
- ecosystem and biodiversity distributions. *PLOS Biology*, 14(3):1–20.
- 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.
- Yesson, C., Brewer, P. W., Sutton, T., Caithness, N., Pahwa, J. S., Burgess, M., Gray, W. A., White, R. J.,
- Jones, A. C., Bisby, F. A., and Culham, A. (2007). How Global Is the Global Biodiversity Information
- ⁸⁰⁷ Facility? *PLoS ONE*, 2(11):e1124.