**Table S1. Description of the dataset according to the ODMAP protocol by Zurell D, Franklin J, König C, et al (2020) A standard protocol for reporting species distribution models. Ecography https://doi.org/10.1111/ecog.04960**

| ODMAP elements | Contents |
|---|---|
| **OVERVIEW** | |
| Authorship | **Authors:** Debojyoti Chakraborty, Norbert Móricz, Ervin Rasztovits, Laura Dobor, Silvio Schueler <br><br> **Contact email:** debojyoti.chakraborty@bfw.gv.at <br> **Title:** <br> **DOI:** |
| Model objective | **SDM Objective:** forecast/transfer <br> **Target output:** probability of occurrence of target tree species |
| Taxon | *Seven tree species of Europe:* <br> *Abies alba, Fagus sylvatica, Larix decidua, Picea abies, Pinus sylvestris, Quercus petraea, Quercus robur* |
| Location | Europe |
| Scale of analysis | **Spatial extent (Lon/ Lat):** <br> Longitude: -32.65000 °E -69.44167 °E <br> Latitude: 30.877982 °N -71.57893 °N <br><br> **Spatial resolution:** 30 arcsec <br><br> **Temporal resolution:** We modeled for historic climate (1961-90) and three future time frames which include averages of (2041-2060, 2061-2080, and 2081-2100). The predictions were done for two Representative Concentration RCP 4.5 and RCP 8.5 |
| Biodiversity data overview | **Observation type:** standardized monitoring <br> **Response data type:** presence/ absence data |
| Type of predictors | Climatic |
| Conceptual model/hypotheses | A large body of scientific studies indicate that climate is one of the major drivers of the distribution of tree species at the continental scale. We exploited this correlation between species' current occurrence and climate to develop SDMs that predict the potential distribution of the target tree species. |
| Assumptions | We assumed that species are at pseudo-equilibrium with the environment. The source of the presence/absence data (Mauri et al. 2017) used in this study is largely from national forest inventories where tree individuals below a certain diameter at |

| | breast height are not recorded. We assume that this data collection procedure did not bias our occurrence data. |
|---|---|
| | Since our occurrence dataset covers the whole current distribution of the target species, which represents both current and likely future climate of Europe, we safely assumed that the species retain their niches across space and time and the current occurrence~ climate correlation remains stable when predicting the models for future climate. |
| SDM algorithms | **Algorithms:** We selected 10 modeling algorithms: GLM (Generalized Linear Models), GAM (Generalized Additive Models), GBM (Generalized Boosted regression Models), CTA (Classification Tree Analysis), ANN (Artificial Neural Networks), SRE (Surface Range Envelop or BIOCLIM), FDA (Flexible Discriminant Analysis), MARS (Multivariate Adaptive Regression Spline), RF (Random Forest for classification and regression), and MAXENT. Tsuruoka. These model algorithms were implemented through an ensemble model platform biomod2 (Thuiller et al. 2016).<br><br>**Model complexity:** The individual models were run using the standard default settings of biomod2, that are designed to balance model complexity and overfitting<br><br>**Ensembles:** The prediction of individual model algorithms were ensembled through biomod2 (Thuiller et al. 2016). |
| Model workflow | The model workflow includes:<br>1. Data cleaning and generation of pseudo absences<br>2. Finding the best climate variables to fit the models<br>2. Model running through biomod2 platform<br>3. Ensemble prediction<br>4. Generation of the maps as gridded 30 arcsec rasters. |
| Software | **Software:** All analyses were conducted using R version 3.3.2 (R Core Team, 2016). Packages used: biomod2 (Thuiller et al. 2016), Random Forest (Breiman 2001),<br><br>Data availability: Presence absence data are available from Mauri et al (2014)<br><br>**Climate data is available from**<br><br>Chakraborty D, Dobor L, A, Hlásny T, Schueler S (2020) High-resolution gridded climate data for Europe based on bias-corrected EURO-CORDEX: the ECLIPS-2.0 dataset [Zenodo: |

| | |
|---|---|
| | |
| **DATA** | |
| Biodiversity data | **Taxon names:** *Abies alba, Fagus sylvatica, Larix decidua, Picea abies, Pinus sylvestris, Quercus petraea, Quercus robur* |
| | **Ecological level:** *Species-level* |
| | **Data source:**<br>Species presence-absence data was obtained from the EU-Forest dataset (Mauri et al. 2017). The dataset harmonizes European tree occurrence from National Forest inventories (NFI), Forest Focus (Hiederer et al. 2011), Biosoil datasets (Houston et al. 2011). A major part of the data arises from the NFI data (96%) while 4% contributed by Forest Focus (Hiederer et al. 2011), Biosoil datasets (Houston et al. 2011). |
| | **Sampling design:** The background data included in the EU-Forest (Mauri et al .2017) varied in their sampling intensity and design. This data was harmonized and aggregated to a spatial resolution of 1 square kilometer, in line with an INSPIREcompliant 1 km× 1 km grid |
| | **Sample size**<br>The dataset includes a total of 1,000,525 occurrence records at a spatial resolution of 1x1km (Mauri et al 2017) |
| | **Data filtering:** Form the EU-Forest dataset we obtained 412,2881 occurrence records about the seven target species. |
| | **Presence-Absence data:**<br>In our case the geographic locations of the target species in the EU-Forest dataset was asumed to be true presences, while the remaining locations of occurrence of other species were asumed to be the absence locations.<br>To ensure that absence locations are not only climatically dissimilar but also geographically distant from the observed presence locations, we developed the so-called pseudo absences according to Senay et al (2013). This is a three-step approach: i) specifying a geographical extent outside the observed presences; ii) environmental profiling of the absences outside this geographic extent, and iii) *k-means* clustering of the environmental profiles and selecting random samples within |

| | |
|---|---|
| | each cluster. In our case, a 2-degree buffer was found to be optimum following Senay et al. (2013). The absence locations outside this geographic extent were classified into 10-15 (depending on species) environmentally dissimilar clusters according to the k-means clustering algorithm.The number of clusters for each species were determined with a plot of total within-cluster sum of square (WSS) and number of clusters.<br><br>The number of pseudoabsence locations was further reduced by randomly selecting a sample of locations defined by the 95% confidence interval from each of the clusters. This approach was used to generate pseudoabsence for all the seven species. |
| Data partitioning | The occurrence dataset for each target species was partitioned by splitting into 75% for model training and 25% for model evaluation. |
| Environmental predictors | **Predictor variables**<br>Environmental predictors were 80 biologically relevant climate variables comprising of annual, seasonal, and monthly variables. From this list of 80 variables, a small subset of potential predictor variables was selected for each target species during the variable selection process.<br><br>**Data sources:**<br>**The spatial resolution of predictor data:** 30 arcsec which is roughly equivalent to 1x1km or lower depending on latitude.<br><br>**The temporal resolution of predictor variable: H**istoric climate (1961-90) and three future time frames which include averages of (2041-2060, 2061-2080 and 2081-2100) for two Representative Concentration RCP 4.5 and RCP 8.5 were used for the SDM predictions.<br><br>**Geographic projection:** WGS 84 (EPSG: 4326) |
| **Model** | |
| Variable selection and Multicollinearity | From the list of potential predictor variables (Table S1), the ones which explain most of the variation in the observed presence and absences of each species were selected with a recursive feature elimination approach (RFE) implemented within the Random forest algorithm (Breiman 2001). Within the RFE approach, the variables were eliminated iteratively, starting from the full set of potential predictors (Table S1), and retaining only those variables that reduce the mean square error over random permutations of the same variable. The variables which were linearly correlated with other variables and had a variance |

| | |
|---|---|
| | inflation factors VIF > 5 as suggested by Booth et al. (1994) were identified and the ones with the lower value according to the Akaike Information Criteria (AIC) (Akaike 1974) were retained for further model development. This subset of uncorrelated climate variables  (Table S2 in Supporting Information) was used as predictor variables for developing the ensemble species distribution models. |
| Model settings | The models were run with the default settings of biomod2 (Thuiller et al 2016) |
| Model estimates | The models estimated median ensemble probability of species occurrence and associated model uncertainty represented by the coefficient of variation. |
| Model ensemble | Predicted probabilities from the individual models for each target species were ensembled as a consensus model which combined the median probability over the selected models with True Skill Statistics threshold (TSS > 0.7) (Allouche et al. 2006; Coetzee et al. 2009). |
| Threshold selection | True Skill Statistics threshold (TSS > 0.7), a commonly used threshold for SDMS (Allouche et al. 2006; Coetzee et al. 2009) was used. |
| **Assessment** | |
| Model performance statistics | For each such model run as well as the final ensemble models for each target species, the model evaluation statistics were recorded. These statistics were true skill statistics (TSS) and area under the relative operating characteristic (ROC), model sensitivity (the ability of the model to predict true presences), and model specificity (the ability of the model to predict the true absences). TSS takes into account both omission and commission errors and ranges also from −1 to +1, not being affected by prevalence as KAPPA (Allouche et al. 2006). TSS values ranging from 0.2 to 0.5 were considered poor, from 0.6 to 0.8 useful, and values larger than 0.8 were good to excellent (e.g. Coetzee et al. 2009). Prediction accuracy is considered to be similar to random for ROC values lower than 0.5; poor, for values in the range 0.5–0.7; fair in the range 0.7–0.9; and excellent for values greater than 0.9 (Pontius and Parmentier 2014). |
| **Prediction** | |
| Prediction output | Predicted probabilities from the individual models and target species were ensembled as a consensus model which combined the median probability over the selected models with True Skill Statistics threshold (TSS > 0.7) (Allouche et al. 2006; Coetzee et al. 2009). The median was chosen because it is known to be less sensitive to outliers than the mean. The estimated ensemble model predictions were presented as GeoTIFF rasters |

| Uncertainty quantification | Model uncertainty was estimated in terms of coefficient of variation (CV) among the predictions of the individual models. The estimated CVs are also presented as GeoTIFF rasters where each cell corresponds to a CV value whereby higher and lower CV values indicate higher and lower uncertainty respectively in the ensemble model. |
| --- | --- |