

Analysing the sensitivity of pollen based land-cover maps to different auxiliary variables

Behnaz Pirzamanbein^{1,2}, Anneli Poska^{3,4}, and Johan Lindström¹

¹Centre for Mathematical Sciences, Lund University, Sweden

²Centre for Environmental and Climate Research, Lund University, Sweden

³Department of Physical Geography and Ecosystems Analysis, Lund University, Sweden

⁴Institute of Geology, Tallinn University of Technology, Estonia

Correspondence to: Behnaz Pirzamanbein (behnaz@maths.lth.se)

Abstract. Realistic depictions of past land cover are needed to investigate prehistoric environmental changes and **anthropogenic impacts** to determine the scale of anthropogenic deforestation on environment and to study long term land cover-climate feedbacks. However, observation based reconstructions of past land cover are rare while commonly used modelled based reconstructions exhibit considerable differences and give diverging results when used by climate modellers to study the effects of past land-cover dynamics on climate. Recently Pirzamanbein et al. (2015, arXiv:1511.06417) developed a statistical interpolation method that produces spatially complete reconstructions of past land cover from pollen assemblage. These reconstructions incorporate a number of auxiliary datasets raising questions regarding both the method's sensitivity to the choice of auxiliary data and the unaffected transmission of observational data.

Here the sensitivity of the method is examined by performing spatial reconstructions for northern Europe during three time periods (1900 CE, 1725 CE and 4000 BCE), based on irregularly distributed pollen based land cover, available for ca 25% of the area, and different auxiliary datasets. The spatially explicit auxiliary datasets considered include the most commonly utilized sources of the past land-cover data — estimates produced by a dynamic vegetation (DVM) and anthropogenic land-cover change (ALCC) models — and modern elevation. Five different auxiliary datasets were considered, including different climate data driving the DVM and different ALCC models. The resulting reconstructions were evaluated using deviance information criteria and cross validation for all the time periods. For the recent time period, 1900 CE, the different land-cover reconstructions were also compared against a present day forest map.

The tests confirm that the developed statistical model provides a robust spatial interpolation tool with low sensitivity to differences in auxiliary data and high capacity to un-distortedly transmit the information provided by sparse pollen based **observations** proxy data. Further, usage of auxiliary data with high spatial detail improves the model performance for the areas with complex topography or where observational data is missing.

1 Introduction

The importance of terrestrial land cover for the global carbon cycle and its impact on the climate system is well recognized (e.g. Claussen et al., 2001; Brovkin et al., 2006; Arneth et al., 2010; Christidis et al., 2013). Many studies have found large

climatic effects associated with changes in land cover. Forecast simulations evaluating the effects of human induced global warming predict a considerable amplification of future climate change, especially for Arctic areas. The amplification is due to a number of biogeophysical and chemical feedbacks brought by the northward advancement of boreal shrub and treeline (Zhang et al., 2013; Richter-Menge et al., 2011; Chapman and Walsh, 2007; Koenigk et al., 2013; Miller and Smith, 2012). The past anthropogenic deforestation of the temperate zone in Europe was lately demonstrated to have an impact on regional climate similar in amplitude to present day climate change (Strandberg et al., 2014). However, studies on the effects of vegetation and land-use changes on past climate and carbon cycle often report considerable differences and uncertainties in their model predictions (de Noblet-Ducoudré et al., 2012; Olofsson, 2013).

One of the reasons for such widely diverging results could be the differences in past land-cover descriptions used by climate modellers. Possible land-cover descriptions range from static present-day land cover (Strandberg et al., 2011), over simulated potential natural land cover from dynamic (or static) vegetation models (DVMs) (e.g. Brovkin et al., 2002; Hickler et al., 2012), to past land-cover scenarios combining DVM derived potential vegetation with estimates of anthropogenic land-cover change (ALCC) (Strandberg et al., 2014; Pongratz et al., 2008; de Noblet-Ducoudré et al., 2012). Differences in input climates, inherent mechanistic and parametrisation differences of DVMs (Prentice et al., 2007; Scheiter et al., 2013), and significant variation in land-use estimates between the existing ALCC scenarios (e.g. Kaplan et al., 2009; Pongratz et al., 2008; Klein Goldewijk et al., 2011; Gaillard et al., 2010) further contribute to the differences in past land-cover descriptions. These differences can lead to largely diverging estimates of past land-cover dynamics even when the most advanced models are used (Strandberg et al., 2014; Pitman et al., 2009).

The palaeoecological ~~observation-proxy~~ based land-cover reconstructions (LCR) recently published by Pirzamanbein et al. (2014, 2015) were designed to overcome the above described problems. And to provide an alternative, ~~observation-proxy~~ based, land-cover description applicable for a range of studies on past vegetation and its interactions with climate, soil and humans. These reconstructions use the pollen based land-cover composition (PbLCC) published by Trondman et al. (2015) as a source of information on past land-cover composition. The PbLCC are point estimates, depicting the land-cover composition of the area surrounding each of the studied sites. To fill the gaps between these observations and to acquire a spatially continuous land-cover reconstruction, spatial interpolation is necessary. Conventional interpolation methods might struggle when handling noisy, spatially heterogeneous data (Heuvelink et al., 1999; De Knecht et al., 2010), but alternative statistical methods for handling spatially structured data exist (e.g. Gelfand et al., 2010; Blangiardo and Cameletti, 2015).

In Pirzamanbein et al. (2015) a statistical model based on Gaussian Markov Random Fields (GMRFs, Lindgren et al., 2011; Rue and Held, 2004) was developed to provide a reliable, computationally effective and freeware based spatial interpolation technique. The resulting statistical model combines PbLCC data with auxiliary datasets; e.g. DVM output, ALCC scenarios, and elevation; to produce reconstructions of past land cover. The auxiliary data is subject to the differences and uncertainties outlined above and the choice of auxiliary data could influence accuracy of the statistical model. ~~This study aims at determining~~
The major objectives of this paper are: 1) to draw attention of climate modelling community to a novel set of spatially explicit proxy based land-cover reconstructions suitable for climate modelling; 2) to present and test the robustness of the model;
statistical spatial interpolation model for proxy based reconstructions developed by Pirzamanbein et al. (2015); 3) to evaluate

its capacity to un-distortedly recover information provided by PbLCC ~~observations~~-proxy data on past vegetation composition; and 4) to analyse the models sensitivity to different auxiliary datasets.

2 Material and methods

The studied area covers temperate, boreal and alpine-arctic biomes of central and northern Europe (45°N to 71°N and 10°W to 30°E). The Pollen based land-cover composition (PbLCC) published in Trondman et al. (2015) consists of proportions of coniferous forest (CF), broadleaved forest (BF) and un-forested land (UF) presented as gridded ($1^\circ \times 1^\circ$) data points placed irregularly across northern-central Europe. Altogether 175 grid cells containing ~~observational~~-proxy data were available for 1900 CE, 181 for 1725 CE, and 196 for the 4000 BCE time-period (Figure 1, column 2).

Four different model derived datasets, depicting past land cover, along with elevation (based on SRTM data) were considered as potential auxiliary datasets. In each case potential natural vegetation (PNV) composition estimated by the dynamic vegetation model (DVM) LPJ-GUESS (Lund-Potsdam-Jena General Ecosystem Simulator; Smith et al., 2001; Sitch et al., 2003) is combined with an ALCC scenario to adjust for human land use (see Pirzamanbein et al., 2014, for more detail):

K-L_{RCA3}: Combines the ALCC scenario KK10 (Kaplan et al., 2009) and the PNV composition from LPJ-GUESS. Climate forcing for the DVM was derived from RCA3 (Rossby Centre Regional Climate Model, Samuelsson et al., 2011) at annual time and $0.44^\circ \times 0.44^\circ$ spatial resolution (Figure 1, column 3),

K-L_{ESM}: Combines the ALCC scenario KK10 and the PNV composition from LPJ-GUESS. For this dataset, the climate forcing for the DVM was derived from the Earth System Model (ESM; Mikolajewicz et al., 2007) at centennial time and $5.6^\circ \times 5.6^\circ$ spatial resolution. To interpolate data into annual time and $0.5^\circ \times 0.5^\circ$ spatial resolution climate data from 1901–1930 CE provided by the Climate Research Unit (CRU) was used (Figure 1, column 4),

H-L_{RCA3}: Combines the ALCC scenario from the History Database of the Global Environment (HYDE; Klein Goldewijk et al., 2011) and the PNV composition from LPJ-GUESS with RCA3 climate forcing (Figure 1, column 5),

H-L_{ESM}: Combines the ALCC scenario from HYDE and the PNV composition from LPJ-GUESS with ESM climate forcing (Figure 1, column 6).

In addition, elevation data used in modelling was obtained from the Shuttle Radar Topography Mission (SRTM_{elev}, Becker et al., 2009) (Figure 1, column 1 row 2).

Finally, a modern forest map based on data from the European Forest Institute (EFI) is used for evaluation of the model's performance for the 1900 CE time period. The EFI forest map (EFI-FM) is based on a combination of satellite data (NOAA-AVHRR) and national forest-inventory statistics from 1990—2005 (Päivinen et al., 2001; Schuck et al., 2002) (Figure 1, column 1 row 1).

All above described sets of auxiliary data were up-scaled to $1^\circ \times 1^\circ$ spatial resolution, matching the pollen based reconstructions, before usage as model input.

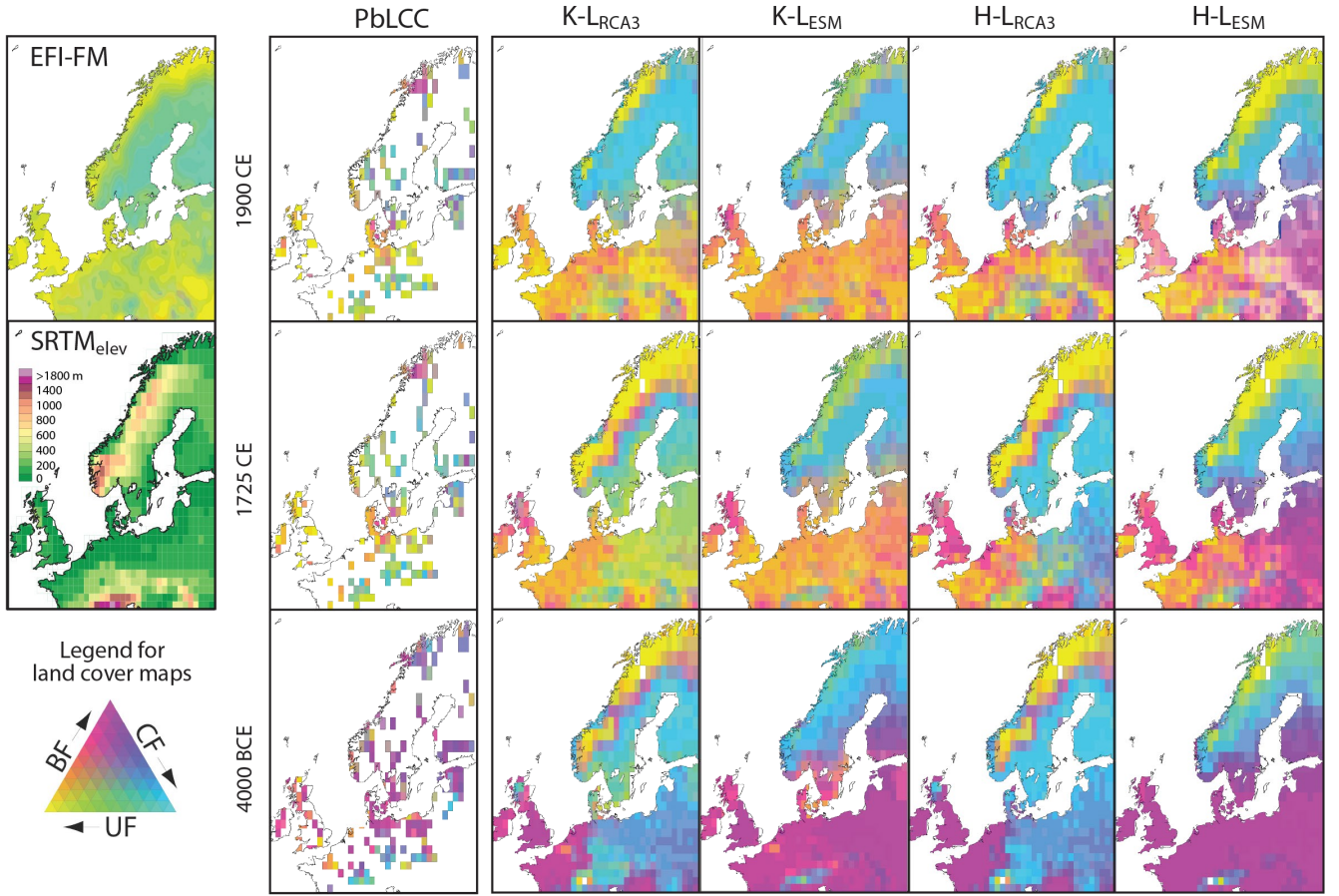


Figure 1. Data used in the modelling. The first column shows (from top to bottom) the EFI forest map, $SRTM_{elev}$, and the colorkey for the land-cover compositions, [coniferous forest \(CF\)](#), [broadleaved forest \(BF\)](#) and [unforested land \(UF\)](#). The remaining columns gives (from left to right) the pollen based land-cover composition (PbLCC, Trondman et al., 2015) and the four model based compositions that could be used as covariates: K- L_{RCA3} , K- L_{ESM} , H- L_{RCA3} , and H- L_{ESM} . Here K/H indicates KK10 (Kaplan et al., 2009) or HYDE (Klein Goldewijk et al., 2011) land use scenarios and L_{RCA3}/L_{ESM} indicates vegetation model driven by climate from the Rossby Centre Regional Climate Model (Samuelsson et al., 2011) or Earth System Model (Mikolajewicz et al., 2007). The three rows representing (from top to bottom) the time periods 1900 CE, 1725 CE, and 4000 BCE.

2.1 Statistical model for land-cover compositions

A Bayesian hierarchical model (Figure 2) is used to model the PbLCC data; [notation used in the model is summarised in Table 1](#). For each component of PbLCC, we assume an underlying compositional vector describing the proportions of land cover: coniferous forest, broadleaved forest and un-forested land. The effect of covariates and spatial structure are incorporated

5 in the underlying compositional vector.

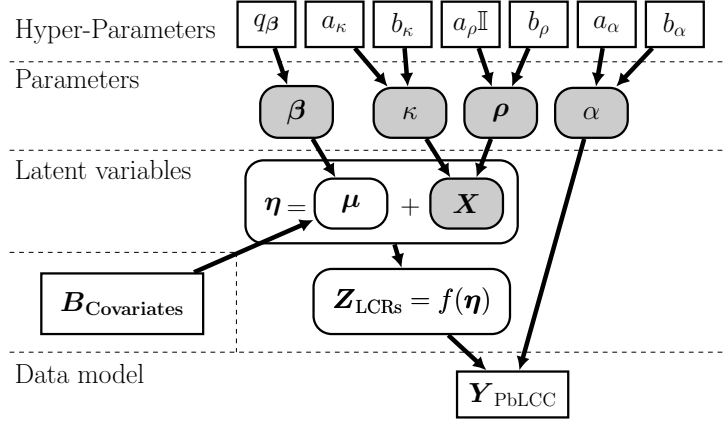


Figure 2. Hierarchical graph describing the conditional dependencies between the model inputs (white rectangle) and parameters (gray rounded rectangle) which need to be estimated. The white rounded rectangle are computed based on the estimations. [The model can be interpreted as an empirical forward model \(direction of arrows\) where parameters affect the latent variables which in turn affect the data. Reconstructions are then obtained by inverting the model \(i.e. computing the posterior\) to obtain the latent variables given the data.](#)

List of notations

α	Concentrated parameter of the Dirichlet distribution (i.e. uncertainty)
β	Regression coefficients
ρ	Covariance matrix that determines the variation between and within fields
κ	Scale parameter controls the range of spatial dependency
μ	Mean structures
X	Spatial dependence residuals
η	Latent Gaussian Markov random fields ($\eta = \mu + X$)
Z_{LCRs}	Vector of proportions or realizations
$B_{\text{Covariates}}$	Matrix of auxiliary datasets
f	Link function (transforms from \mathbb{R}^2 to proportions $(0,1)^3$; $Z = f(\eta)$)
Y_{PbLCC}	Observations, as proportions.

Table 1. [Summary of notation used in the Bayesian hierarchical model.](#)

To account for observational uncertainty in the compositions, the PbLCC are modelled as draws from a Dirichlet distribution given concentrated parameter α (controlling the uncertainty) and the vector of proportions \mathbf{Z} ,

$$\mathbf{Y}_{\text{PbLCC}}|\mathbf{Z}, \alpha \sim \text{Dir}(\alpha\mathbf{Z}) \quad \alpha > 0, \mathbf{Z}_k \in (0, 1), \sum_k \mathbf{Z}_k = 1.$$

To account for the spatial dependence in the proportions, \mathbf{Z} is modelled as a transformation, f , of a latent GMRF, $\boldsymbol{\eta}$:

$$\mathbf{Z} = f(\boldsymbol{\eta}) \quad f: \mathbb{R}^2 \rightarrow (0, 1)^3$$

$$\mathbf{Z}_k = \begin{cases} \frac{\exp(\boldsymbol{\eta}_k)}{1 + \sum_{i=1}^2 \exp(\boldsymbol{\eta}_i)} & \text{for } k = 1, 2 \\ \frac{1}{1 + \sum_{i=1}^2 \exp(\boldsymbol{\eta}_i)} & \text{for } k = 3 \end{cases}.$$

The inverse of f is called the additive log-ratio transformation (alr, Aitchison, 1986), i.e. $\boldsymbol{\eta}_k = \log(\mathbf{Z}_k/\mathbf{Z}_3)$, $k = 1, 2$. The alr transformation is the multivariate extension of a logit transformation.

The latent field is modelled with a mean structure $\boldsymbol{\mu}$ and a spatially dependent residual \mathbf{X} ,

$$\boldsymbol{\eta} = \begin{bmatrix} \boldsymbol{\eta}_1 \\ \boldsymbol{\eta}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} + \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} = \mathbf{X} + \boldsymbol{\mu}.$$

Here \mathbf{X} is GMRF with a separable covariance structure;

$$\mathbf{X}|\kappa, \boldsymbol{\rho} \sim \text{N}(0, \boldsymbol{\rho}^{-1} \otimes \mathbf{Q}(\kappa))$$

where $\mathbf{Q}(\kappa)$ is the precision matrix of a spatially dependent GMRF (Lindgren et al., 2011), κ is the scale parameter which controls the range of spatial dependency and $\boldsymbol{\rho}$ controls the variation within and between the fields \mathbf{X}_k (see Pirzamanbein et al., 2015, for details).

The mean structure is modelled as a linear regression $\boldsymbol{\mu} = \mathbf{B}\boldsymbol{\beta}$, i.e. a combination of covariates \mathbf{B} and regression coefficients $\boldsymbol{\beta}$. The main focus of this paper is to evaluate the model sensitivity to the choice of covariates (e.g. the auxiliary datasets). The PbLCC is modelled based on six different sets of covariates (Figure 1): 1) Intercept, 2) SRTM_{elev}, 3) K-L_{ESM}, 4) K-L_{RCA3}, 5) H-L_{ESM}, and 6) H-L_{RCA3}. Table 2 shows the different models and the corresponding covariates included in the model.

The model description is completed by specifying prior distributions for the model parameters. Wide but proper priors are assigned for $\alpha, \kappa, \boldsymbol{\rho}$ and $\boldsymbol{\beta}$. Specifically, a Gamma prior is chosen for the uncertainty and scale parameters, α and κ , i.e. $\Gamma(1.5, 0.1)$ and $\Gamma(1.5, 0.1)$. A Gaussian prior for the regression parameters $\boldsymbol{\beta}$, with zero expectation and small precision $q_{\boldsymbol{\beta}} = 10^{-3}$. The $\boldsymbol{\rho}$ is assigned an inverse wishart prior, $IW(\mathbb{I}, 10)$, where \mathbb{I} is a 2×2 identity matrix.

2.2 Inference and associated uncertainties

The Markov Chain Monte Carlo (MCMC) method (Brooks et al., 2011) is used to estimate the parameters and to reconstruct the land-cover composition, \mathbf{Z}_{LCRs} , with 100 000 MCMC samples and a burn-in sample size of 10 000. Details of the MCMC implementation can be found in Pirzamanbein et al. (2015).

Model	Covariates					
	Intercept	SRTM _{elev}	K-L _{ESM}	K-L _{RCA3}	H-L _{ESM}	H-L _{RCA3}
Constant	x					
Elevation	x	x				
K-L _{ESM}	x	x	x			
K-L _{RCA3}	x	x		x		
H-L _{ESM}	x	x			x	
H-L _{RCA3}	x	x				x

Table 2. Six different models and corresponding covariates. SRTM_{elev} is elevation (Becker et al., 2009), K/H indicates KK10 (Kaplan et al., 2009) or HYDE (Klein Goldewijk et al., 2011) land use scenarios and L_{RCA3}/L_{ESM} indicates vegetation model driven by climate from the Rossby Centre Regional Climate Model (Samuelsson et al., 2011) or Earth System Model (Mikolajewicz et al., 2007).

In each MCMC iteration, the samples of η are obtained by adding the spatial dependency field \mathbf{X} and the effect of covariates through $\mathbf{B}\beta$. Applying the alr transformation to the η samples, MCMC samples for \mathbf{Z} are obtained. The land-cover reconstruction is then computed by averaging the MCMC samples, giving $\mathbf{Z}_{\text{LCRs}} = \mathbf{E}(\mathbf{Z}|\mathbf{Y}_{\text{PbLCC}})$.

The uncertainties of the land-cover reconstruction, $\mathbf{V}(\mathbf{Z}|\mathbf{Y}_{\text{PbLCC}})$, are assessed by constructing predictive regions (PR) using the MCMC samples at each location. The predictive regions are constructed to represent the uncertainty associated with the reconstructions; including uncertainties in both model parameters (α , κ , ρ , and β) and underlying fields (\mathbf{Z}). The predictive regions are constructed by first creating elliptical 95%-predictive regions (i.e. containing 95% of the MCMC samples) in \mathbb{R}^2 (Figure 3 left plot) before transforming these to ternary predictive region in $(0,1)^3$ (Figure 3 right plot) (see Pirzamanbein et al., 2015, for details). In order to compare the uncertainties of different model land-cover reconstructions, we report the fraction of the unit triangle covered by the ternary PR. This is done by distributing points in the ternary diagram and computing the fraction as the number of points laying inside the PR divided by total number of points in the ternary triangle.

2.3 Testing the model performance

To evaluate the model performance, we compared the land-cover reconstructions from different models for the 1900 CE time period with the European Forest Institute forest map (EFI-FM) by computing the average compositional distances (ACD). The compositional distances between two different compositions, \mathbf{U} and \mathbf{V} , are computed as (Aitchison et al., 2000)

$$\Delta(\mathbf{U} - \mathbf{V}) = \Delta(\mathbf{u} - \mathbf{v}) = \left((\mathbf{u} - \mathbf{v})^\top \mathbf{H}^{-1} (\mathbf{u} - \mathbf{v}) \right)^{1/2}$$

where $\mathbf{u} = \text{alr}(\mathbf{U})$, $\mathbf{v} = \text{alr}(\mathbf{V})$ and \mathbf{H} is a 2×2 matrix, neutralizing the choice of denominator in the alr transformation, with elements $H_{ij} = 2$ if $i = j$, and $H_{ij} = 1$ if $i \neq j$. These distances are then averaged over all locations. This measure is similar to root mean square error in \mathbb{R}^2 but it accounts for compositional properties, i.e. each component of the compositions is between $(0,1)$ and sum of all the components is 1.

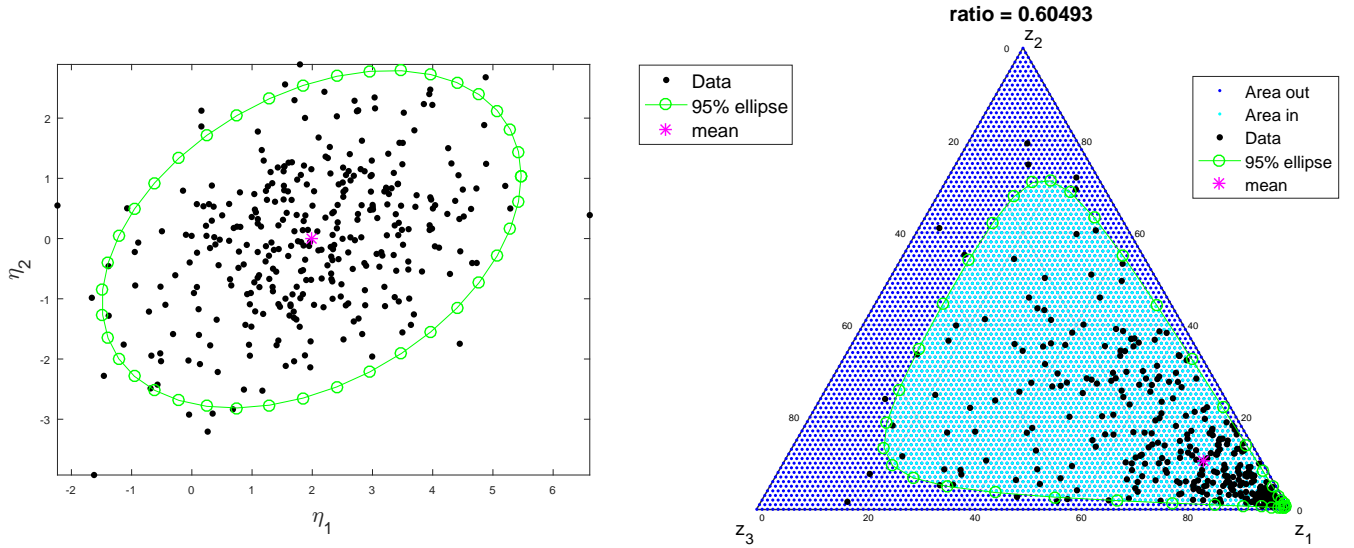


Figure 3. The left plot shows the 95% elliptical predictive region in \mathbb{R}^2 . The right ternary diagram shows the transformed 95% predictive region together with the corresponding fraction, 60%, compared to the whole triangle.

Since no independent observational data exists for the 1725 CE and 4000 BCE time periods, we applied a 6-fold cross-validation scheme (Friedman et al., 2001, Ch. 7.10) for all six models and three time periods. The PbLCC data were divided into 6 randomly selected groups and, in each round, the distance between the left out data, $\mathbf{Y}_{\text{PbLCC},l}$, and predictions for group l given a model fitted to the rest of the data, $E(\mathbf{Z}_l | \mathbf{Y}_{\text{PbLCC},k} \text{ } k \notin l)$, were computed.

- 5 To compare the predictive performance of the models, the Deviance Information Criteria (DIC; see Gelman et al., 2014, Ch. 7.2) is also computed for all models and time periods.

3 Results and discussion

Fossil pollen is a well-recognized information source of vegetation dynamics and generally accepted as the best observational data source on past land-cover composition and structure-environmental conditions (Trondman et al., 2015).

- 10 Today, central and northern Europe have, at the subcontinental spatial scale, the highest density of palynologically investigated sites on Earth. However, the collection of pollen data is very time consuming and, cannot be performed everywhere and hardly applicable by climate modellers as input in their original format and therefore heavily underused. The lack of spatially explicit proxy based land cover data directly usable in climate models has been hampering the correct representation of past climate-land cover relationship. This considerably restricts the applicability of the dataset; especially if continuous descriptions of vegetation cover are required. The purpose of the statistical model described in section 2.1 is to combine the observed PbLCC with the spatial structure in the auxiliary data to produce spatially complete reconstructions of past land-cover.
- 15

Regrettably, the available auxiliary datasets (Table 2) exhibit large variation in the extent of coniferous and broadleaved forests, and un-forested areas for all of the studied time periods (Figure -1). These substantial differences illustrate large deviations between model based estimates of the past land-cover composition due to differences in climate forcing and/or applied ALCC scenarios. Differences in climate model outputs (e.g. Harrison et al., 2014; Gladstone et al., 2005) and ALCC model estimates (Gaillard et al., 2010) have been recognized in earlier comparison studies and syntheses. The effect of the differences in input climate forcing and ALCC scenario on DVM estimated land-cover composition presented here are especially pronounced for central and western Europe, and for elevated areas in northern Scandinavia and the Alps (Figure -1). In general the KK10 ALCC scenario produces larger un-forested areas, notably in western Europe, compared to the HYDE scenario. Compared to the ES climate forcing; the RCA3 forcing results in higher proportions of coniferous forest, especially for central, northern and eastern Europe. The described differences are clearly recognizable for all the considered time periods and are generally larger between time periods than within each time period.

Usage of the above described, solely model based land cover to assess the impact of past anthropogenic changes on climate and terrestrial nutrient cycles has been shown to lead to largely diverging results (e.g. Strandberg et al., 2011). The importance of reliable land-cover representations for studying the biogeophysical impacts of anthropogenic land-cover change on climate is well recognized by the climate modelling community (e.g. Strandberg et al., 2011; Pitman et al., 2009; de Noblet-Ducoudré et al., 2012). Introducing the pollen based land cover reconstructions into climate models would facilitate the mechanistic studies on past climate-land cover relationships. The purpose of the statistical model described in Section 2.1 is to combine the observed PbLCC with the spatial structure in the auxiliary data to produce spatially complete maps of past land-cover that can be used directly in climate models.

To assess the impact of the different auxiliary datasets, the statistical model was used to create a set of ~~observation-proxy~~ based reconstructions of past land cover for central and northern Europe during three time periods (1900 CE, 1725 CE and 4000 BCE; see Figures 4–6). Each of the reconstructions were based on the irregularly distributed observed pollen data (PbLCC), available for ca 25% of the area, together with one of the six models described in Table 2; each model uses a different combination of auxiliary data (Figure 1). The spatial dependence in the PbLCC data was modelled using GMRFs, along with additional spatial structure inferred from the auxiliary datasets (Pirzamanbein et al., 2015).

To illustrate the structure of the statistical model, step by step advancement from auxiliary data (model derived land cover) to final statistical estimates, for 1725 CE, are given in Figures 7 and 8. Figure 7 shows, for two locations, how the large differences in K-LRCA3 and K-LESM are reduced by scaling with the regression coefficients, β ; capturing the empirical relationship between covariates and PbLCC data. The two land-cover estimates are then further subject to similar adjustments due to intercept and $SRTM_{elev}$, and finally similar spatial dependent effects. The corresponding decomposition of contributions to the continental map for 1725 CE, from $SRTM_{elev}$, K-LESM, and the spatial effects are given in Figure 8.

The final land-cover reconstructions achieved by fitting the models to the observed PbLCC are very similar, regardless of the auxiliary datasets used.

At first the similarity among the reconstructions might seem contradictory, but recall that the model allows for, and estimates, different weighting (the regression coefficients, β :s) for each of the auxiliary datasets. Thus, the resulting reconstruction do not

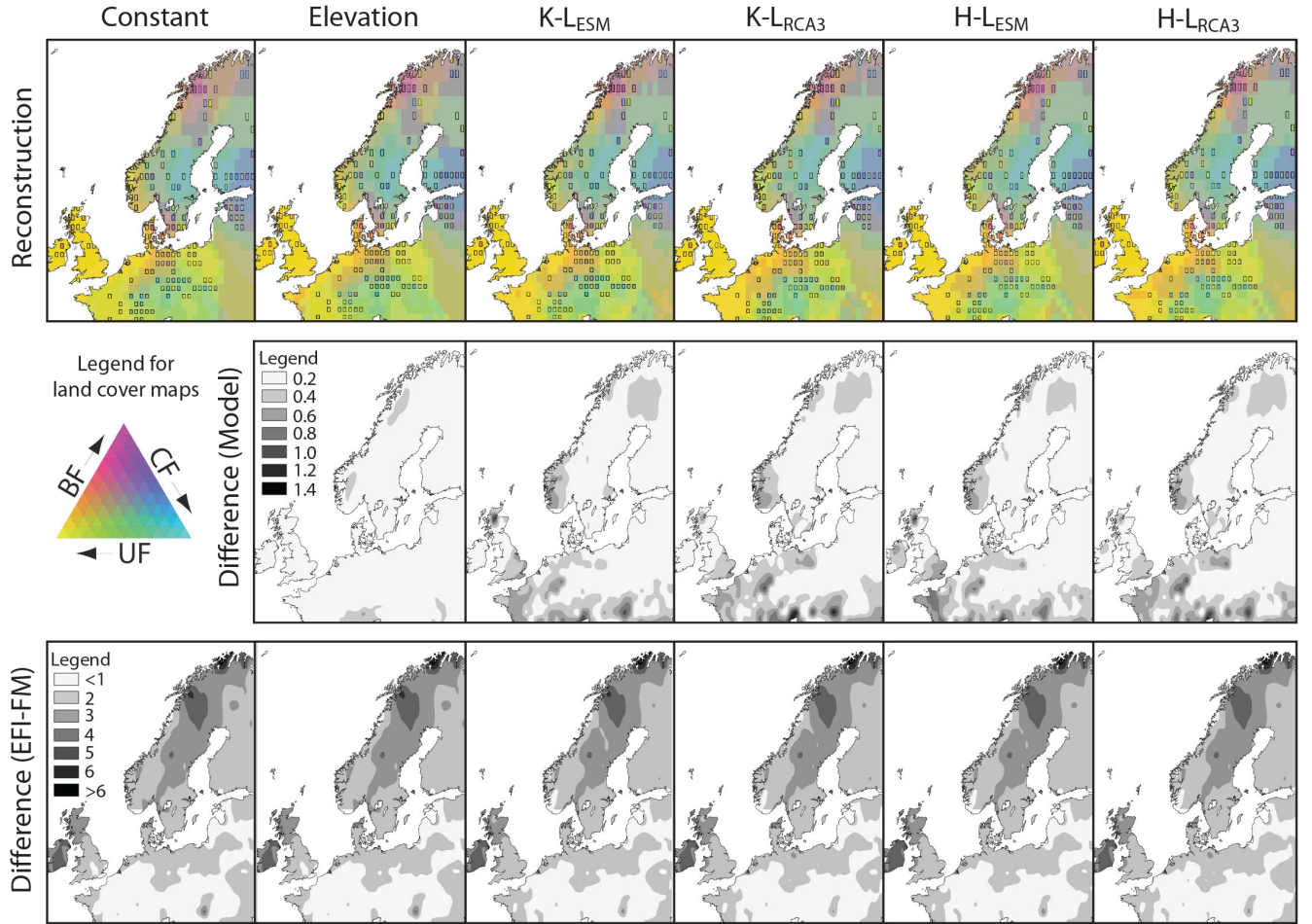


Figure 4. Land-cover reconstructions using pollen based land-cover compositions (PbLCC) for the 1900 CE time periods (top row). The reconstructions are based on six different models (see Table 2) with different auxiliary datasets. Locations and compositional values of the available PbLCC data are given by the black rectangles. Middle row shows the compositional distances between each model and the Constant model. Bottom row shows the compositional distances between each model and the EFI-FM.

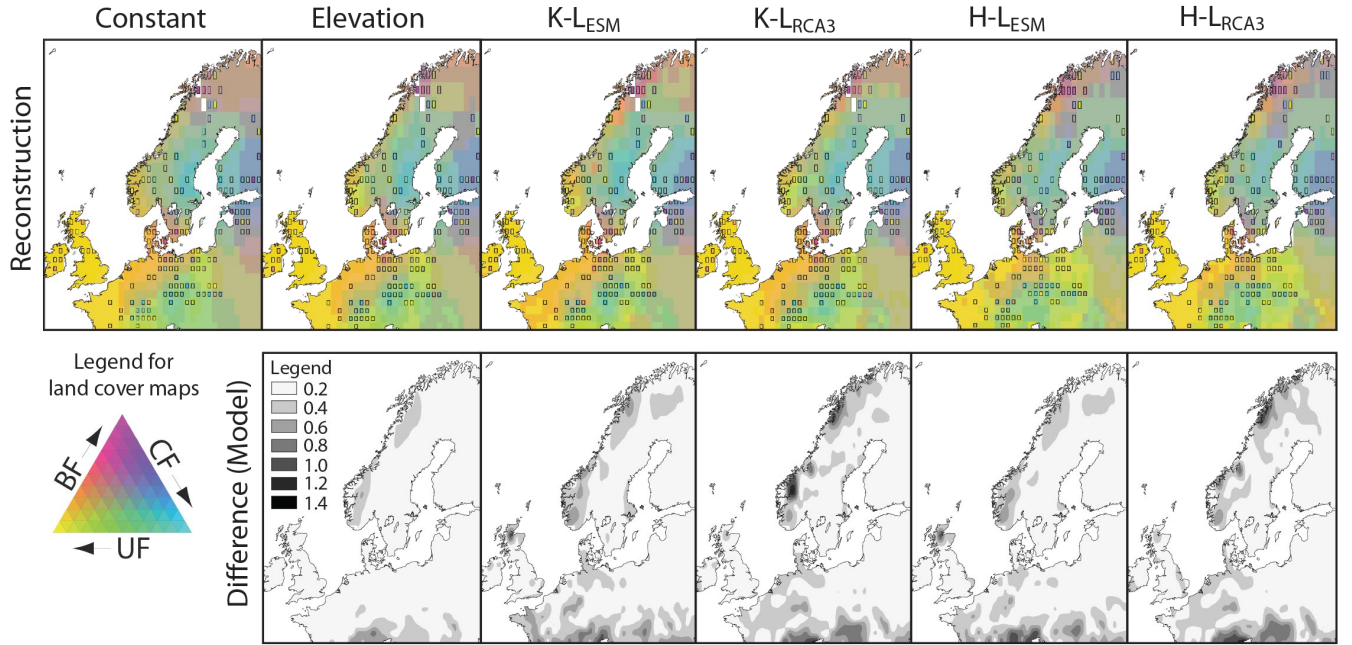


Figure 5. Land-cover reconstructions using local estimates of pollen based land-cover compositions (PbLCC) for the 1725 CE time period (top row). The reconstructions are based on six different models (see Table 2) with different auxiliary datasets. Locations and compositional values of the available PbLCC data are given by the black rectangles. Bottom row shows the compositional distances between each model and the Constant model.

rely on the absolute values in the auxiliary datasets, only their spatial patterns; Table 3 illustrates the substantial discrepancies in the estimated coefficients, β .

<u>1725 CE</u>	<u>Constant</u>	<u>Elevation</u>	<u>K-ESM</u>	<u>K-LRCA3</u>	<u>H-ESM</u>	<u>H-LRCA3</u>
<u>$\beta_{\text{intercept},1}$</u>	-1.39	-1.41	-0.70	-1.20	-0.77	-1.25
<u>$\beta_{\text{intercept},2}$</u>	-1.01	-1.01	-0.73	-1.05	-0.69	-1.01
<u>$\beta_{\text{SRIM}_{\text{elev}},1}$</u>		0.08	0.17	0.13	0.17	0.12
<u>$\beta_{\text{SRIM}_{\text{elev}},1}$</u>		-0.16	0.07	-0.01	0.08	-0.02
<u>$\beta_{\bullet 1,1}$</u>			0.01	-0.07	0.02	0.02
<u>$\beta_{\bullet 1,2}$</u>			-0.02	-0.01	0.01	0.13
<u>$\beta_{\bullet 2,1}$</u>			-0.01	0.10	0.03	0.15
<u>$\beta_{\bullet 2,2}$</u>			-0.19	-0.17	-0.18	-0.12

Table 3. The estimated β coefficient for the six models (see Table 2) for 1725 CE time period. The last four rows represents the coefficients for the modelled derived data set used in that model represented by $\beta_{\bullet i,j}$.

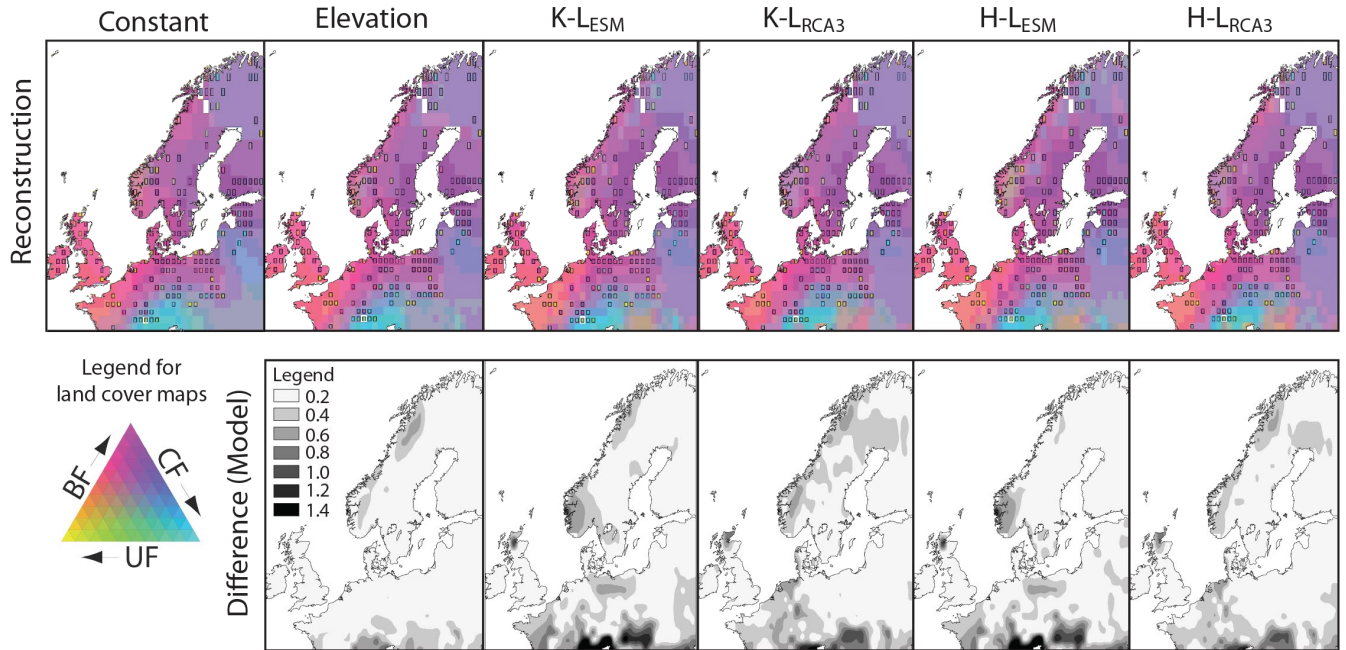


Figure 6. Land-cover reconstructions using local estimates of pollen based land-cover compositions (PbLCC) for the 4000 BCE time period (top row). The reconstructions are based on six different models (see Table 2) with different auxiliary datasets. Locations and compositional values of the available PbLCC data are given by the black rectangles. Bottom row shows the compositional distances between each model and the Constant model.

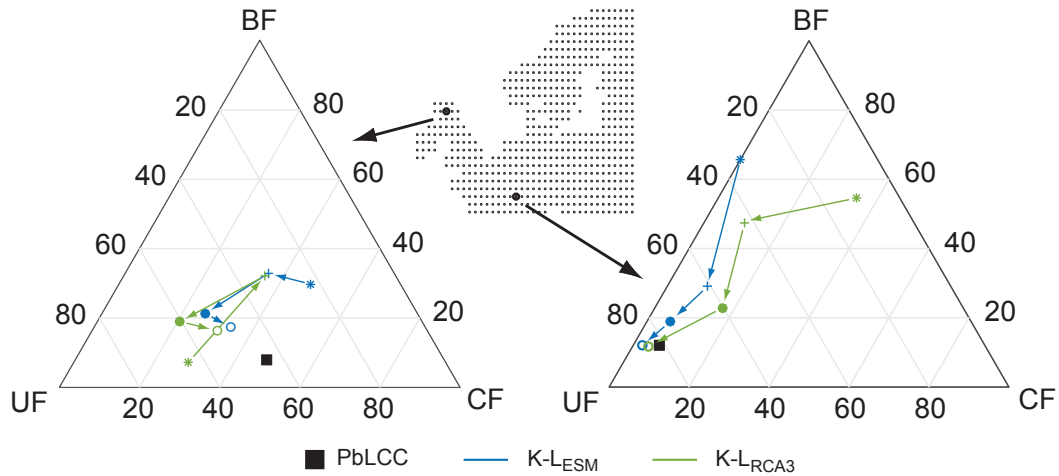


Figure 7. Advancement of the model for two locations at 1725 CE. Starting from the value of the K-LRCA3 and K-LESM covariates (*), the cumulative effects of regression coefficients, β , (+); the intercept and SRTM_{elev} covariates (•); and, finally, the spatial dependency structures (◦), are illustrated. With the final points (◦) corresponding to the land-cover reconstructions, Z_{LCRs} , and ■ marking the observed pollen based land-cover composition.

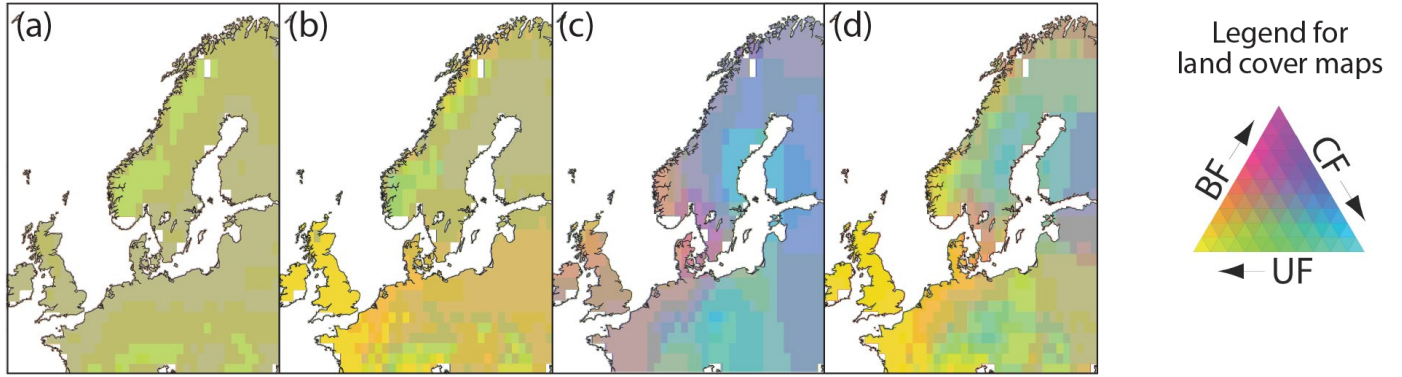


Figure 8. Advancement of K-L_{ESM} models for the 1725 CE time period: (a) shows the effect of intercept and SRTM_{elev}, (b) shows the mean structure, μ , including all the covariates, (c) shows the spatial dependency structure and finally (d) shows the resulting land-cover reconstructions, Z_{LCRs} , obtained by adding (b) and (c).

Although the land-cover reconstructions produced by different models are very similar, model performance for elevated areas and for the areas with low observational data coverage (e.g. eastern and south-eastern Europe) is improved by including covariates that exhibit distinct spatial structures for the given areas (Figures 4–6).

The resulting land-cover reconstructions exhibit considerably higher similarity with the PbLCC data than the auxiliary land-cover datasets for all tested models and time periods (Figures 4–6). The predictive regions indicate the capability of all the models in capturing the PbLCC data and shows similar reconstruction uncertainties (Figure 9). Analogous to the reconstructions the uncertainty regions (the transformed ellipses are described in Figure 3) are very similar in both size and shape irrespective of the auxiliary dataset used. Further the PbLCC data almost always falls within the uncertainty region illustrating that the reconstructions are consistent with the data.

~~The~~ Although a temporal misalignment exists between the PbLCC data for the 1900 CE time period (based on pollen data from 1850 to the present) and the EFI-FM (inventory and satellite data from 1990-2005); EFI-FM provides the best complete and consistent land cover map of Europe for present time, making it a reasonable choice for the comparison. The main differences between the EFI-FM and the PbLCC data for the 1900 CE time period are: 1) lower abundance of broadleaved forests for most of Europe, 2) higher abundance of coniferous forest in Sweden and Finland, and 3) higher abundance of unforested land in North Norway in the EFI-FM data than in the PbLCC data (Pirzamanbein et al., 2015). The average compositional distances computed between the land-cover reconstructions and the EFI-FM for 1900 CE show practically identical (1.47 to 1.48) distances between all six reconstructions and the EFI-FM, and small differences among the six presented models (Table 4). Note that the compositional differences for each grid cells are shown in Figures 4–6. Neither the DIC results (Table 5) nor the 6-fold cross validation results (Table 6) show any advantage among the six tested models for the different time periods. Implying there is no clear preference among the models. These results clearly show that the developed statistical

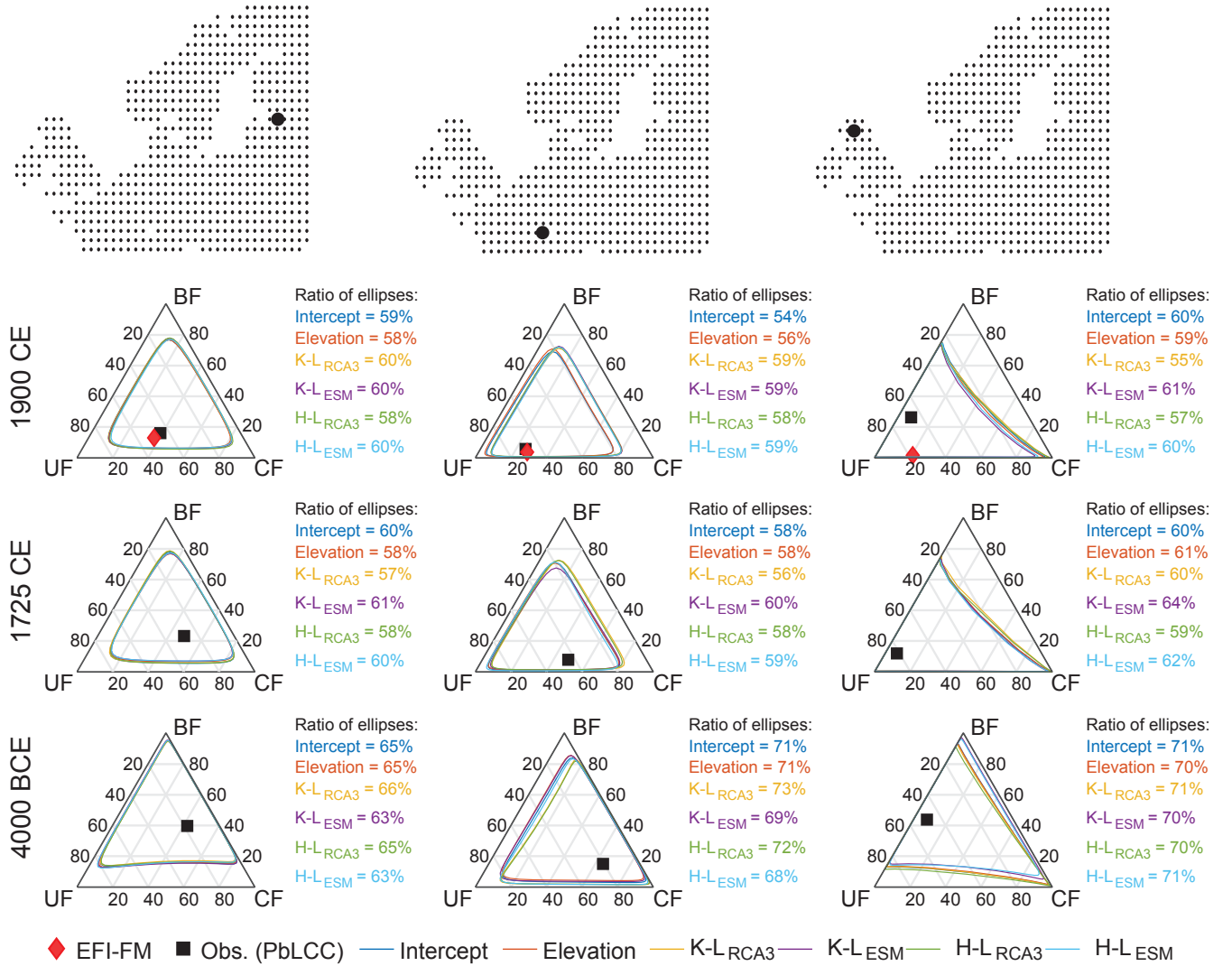


Figure 9. The prediction regions and fraction of the ternary triangle covered by these regions are presented for three locations, the six models (see Table 2), and the 1900 CE, 1725 CE and 4000 BCE time periods. [Construction and interpretation of the prediction regions are described in Section 2.2 and Figure 3.](#)

interpolation model is robust to the choice of covariates. The model is suitable for reconstructing spatially continuous maps of past land cover from scattered and irregularly spaced pollen based ~~observations~~proxy data.

ACD						
Model	1900 CE					
	EFI-FM	Elevation	K-L _{ESM}	K-L _{RCA3}	H-L _{ESM}	H-L _{RCA3}
Constant	1.47	0.06	0.19	0.17	0.19	0.18
Elevation	1.48		0.18	0.16	0.18	0.17
K-L _{ESM}	1.47			0.08	0.10	0.07
K-L _{RCA3}	1.47				0.06	0.11
H-L _{ESM}	1.47					0.08
H-L _{RCA3}	1.47					
1725 CE						
Constant		0.11	0.19	0.16	0.19	0.18
Elevation			0.12	0.14	0.14	0.16
K-L _{ESM}				0.15	0.16	0.07
K-L _{RCA3}					0.08	0.18
H-L _{ESM}						0.17
4000 BCE						
Constant		0.12	0.19	0.21	0.21	0.23
Elevation			0.12	0.19	0.16	0.21
K-L _{ESM}				0.19	0.21	0.07
K-L _{RCA3}					0.07	0.19
H-L _{ESM}						0.21

Table 4. The average compositional distances among the six models (see Table 2) fitted to the data for each of the three time periods.

DIC	1900 CE	1725 CE	4000 BCE
Constant	-562	-656	-591
Elevation	-557	-668	-590
K-L _{ESM}	-551	-654	-601
K-L _{RCA3}	-559	-673	-588
H-L _{ESM}	-554	-654	-607
H-L _{RCA3}	-559	-672	-594

Table 5. Deviance information criteria (DIC) for each of the six models (see Table 2) and three time periods.

ACD	1900 CE	1725 CE	4000 BCE
Constant	0.98	1.13	1.19
Elevation	0.98	1.11	1.20
K-L _{ESM}	0.99	1.12	1.18
K-L _{RCA3}	0.99	1.13	1.18
H-L _{ESM}	1.00	1.12	1.17
H-L _{RCA3}	0.97	0.97	1.17

Table 6. Average compositional distances from 6-fold cross-validations for each of the six models (see Table 2), and three time periods.

4 Conclusions

The ~~performance of the statistical model, explained in Section 2, to reconstruct the pollen-based observations was tested in order to analyse its sensitivity to auxiliary datasets. The observations are~~ statistical model and Bayesian interpolation method presented here has been specially designed for handling palaeo-proxy records and, dependent on proxy data availability, is globally applicable. The ability of the model to create pollen based land cover reconstructions at sub-continental scale was illustrated using an example application on two Holocene time periods frequently assessed by climate researchers: the Little Ice Age (1725 CE) and the Holocene Thermal Maximum (4000 BCE); as well as a third time period covering the recent past (1900 CE) which was used for model validation. The model combines irregularly distributed, pollen based estimates of land cover representing 25% of the study area. ~~The model combines these observations,~~ with auxiliary data and estimates spatial dependencies to produce land-cover maps. The resulting maps capture important features in the pollen proxy data and are reasonably insensitive to the use of different auxiliary datasets.

The considered auxiliary datasets were compiled using most commonly utilized sources of the past land-cover data (estimates produced by a dynamic vegetation model and anthropogenic land-cover changes scenarios). These datasets exhibit considerable model and/or input dependant differences in their recreation of the past land cover. Emphasizing the need for the independent, ~~observation-proxy~~ based past land-cover maps created in this paper.

The ~~model~~-results also indicate that the model has a very good performance and will be very useful for large-scale, continental reconstructions of past land cover. The spatial model tested in this paper can provide an important tool to generate the regional to global scale land-cover maps based on proxy data. Such, pollen based past land cover reconstructions with global coverage are currently produced by the PAGES (Past Global changes) LandCover6k initiative ¹ for most of globe.

5 The model's sensitivity to usage of different auxiliary datasets was validated by calculating deviance information criteria (DIC) and using cross validation for all the time periods. For the recent time period, 1900 CE, the land-cover reconstructions from the different models were also compared against a present day forest map. The evaluation indicates that the applied statistical model is robust. The model estimates the empirical relationship between auxiliary data and pollen based ~~observations~~proxy data, therefore, it only considers features in the auxiliary data that are consistent with the ~~observed~~-proxy data. The covariates with detailed spatial information improves the interpolation results for areas with low ~~observational~~-proxy data coverage, however the overall performance remains unchanged.

Therefore, the model can provide reliable results using a variety of land cover data sets that capture important spatial patterns from vegetation models and past human land use, absent good covariates elevation can be used as the only auxiliary dataset. An important feature of the suggested model is the estimation of different weights for each of the auxiliary datasets (see Table 3).
15 thus capturing the spatial patterns and not the absolute values in the auxiliary datasets. Our validations indicate that auxiliary datasets obtained using different climatic drivers produce very similar reconstructions, which are all close to the pollen based proxy data.

This modelling approach has strong ability to produce empirically based land-cover reconstructions for climate modelling purposes. Such reconstructions are necessary to evaluate the anthropogenic land-cover change scenarios currently used in the climate modelling and to study interactions between land cover and climate in the past with greater reliability.

5 Data availability

The database containing the reconstructions of coniferous forest, broadleaved forest and un-forested land, three fractions of land cover, for the three time-periods presented in this paper, along with reconstructions for 1425 CE and 1000 BCE using only the K-L_{ESM} are available for download from <https://behnaz.pirzamanbin.name/phd/>.

25 *Acknowledgements.* The research presented in this paper is a contribution to the two Swedish strategic research areas Biodiversity and Ecosystems in a Changing Climate (BECC), and Modelling the Regional and Global Earth system (MERGE).

Lindström has been funded by Swedish Research Council (Vetenskapsrådet) grant no 2012-5983.

The authors would like to acknowledge Marie-José Gailard for her efforts in providing the pollen based land-cover ~~observations~~-proxy data and thank her for valuable comments on this manuscript.

¹ www.pastglobalchanges.org/ini/wg/landcover6k/intro

References

- Aitchison, J.: The statistical analysis of compositional data, Chapman & Hall, Ltd., 1986.
- Aitchison, J., Barceló-Vidal, C., Martín-Fernández, J., and Pawłowsky-Glahn, V.: Logratio analysis and compositional distance, *Math. Geol.*, 32, 271–275, 2000.
- 5 Arneth, A., Harrison, S. P., Zaehle, S., Tsigaridis, K., Menon, S., Bartlein, P. J., Feichter, J., Korhola, A., Kulmala, M., O'donnell, D., et al.: Terrestrial biogeochemical feedbacks in the climate system, *Nat. Geosci.*, 3, 525–532, 2010.
- Becker, J., Sandwell, D., Smith, W., Braud, J., Binder, B., Depner, J., Fabre, D., Factor, J., Ingalls, S., Kim, S., et al.: Global bathymetry and elevation data at 30 arc seconds resolution: SRTM30_PLUS, *Marine Geodesy*, 32, 355–371, 2009.
- Blangiardo, M. and Cameletti, M.: Spatial and Spatio-temporal Bayesian Models with R-INLA, Wiley, 2015.
- 10 Brooks, S., Gelman, A., Jones, G. L., and Meng, X.-L.: Handbook of markov chain monte carlo, Chapman and Hall/CRC, 2011.
- Brovkin, V., Bendtsen, J., Claussen, M., Ganopolski, A., Kubatzki, C., Petoukhov, V., and Andreev, A.: Carbon cycle, vegetation, and climate dynamics in the Holocene: Experiments with the CLIMBER-2 model, *Global. Biogeochem. Cy.*, 16, 2002.
- Brovkin, V., Claussen, M., Driesschaert, E., Fichetef, T., Kicklighter, D., Loutre, M., Matthews, H., Ramankutty, N., Schaeffer, M., and Sokolov, A.: Biogeophysical effects of historical land cover changes simulated by six Earth system models of intermediate complexity, *Clim. Dynam.*, 26, 587–600, 2006.
- 15 Chapman, W. L. and Walsh, J. E.: Simulations of Arctic temperature and pressure by global coupled models, *J. Climate*, 20, 609–632, 2007.
- Christidis, N., Stott, P. A., Hegerl, G. C., and Betts, R. A.: The role of land use change in the recent warming of daily extreme temperatures, *Geophys. Res. Lett.*, 40, 589–594, 2013.
- Claussen, M., Brovkin, V., and Ganopolski, A.: Biogeophysical versus biogeochemical feedbacks of large-scale land cover change, *Geophys. Res. Lett.*, 28, 1011–1014, 2001.
- 20 De Knegt, H., van Langevelde, F. v., Coughenour, M., Skidmore, A., De Boer, W., Heitkönig, I., Knox, N., Slotow, R., Van der Waal, C., and Prins, H.: Spatial autocorrelation and the scaling of species–environment relationships, *Ecology*, 91, 2455–2465, 2010.
- de Noblet-Ducoudré, N., Boisier, J.-P., Pitman, A., Bonan, G., Brovkin, V., Cruz, F., Delire, C., Gayler, V., Van den Hurk, B., Lawrence, P., et al.: Determining robust impacts of land-use-induced land cover changes on surface climate over North America and Eurasia: results from the first set of LUCID experiments, *J. Climate*, 25, 3261–3281, 2012.
- 25 Friedman, J., Hastie, T., and Tibshirani, R.: The elements of statistical learning, vol. 1, Springer series in statistics Springer, Berlin, 2001.
- Gaillard, M.-J., Sugita, S., Mazier, F., Trondman, A.-K., Brostrom, A., Hickler, T., Kaplan, J. O., Kjellström, E., Kokfelt, U., Kuneš, P., , Lemmen, C., Miller, P., Olofsson, J., Poska, A., Rundgren, M., Smith, B., Strandberg, G., Fyfe, R., Nielsen, A., Alenius, T., Balakauskas, L., Barnekov, L., Birks, H., Bjune, A., Björkman, L., Giesecke, T., Hjelle, K., Kalnina, L., Kangur, M., van der Knaap, W., Koff, T., Lagerås, P., Latalowa, M., Leydet, M., Lechterbeck, J., Lindbladh, M., Odgaard, B., Peglar, S., Segerström, U., von Stedingk, H., and Seppä, H.: Holocene land-cover reconstructions for studies on land cover-climate feedbacks, *Clim. Past.*, 6, 483–499, 2010.
- 30 Gelfand, A., Diggle, P. J., Guttorp, P., and Fuentes, M.: Handbook of spatial statistics, CRC Press, 2010.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B.: Bayesian data analysis, CRC press, 2014.
- Gladstone, R., Ross, I., Valdes, P., Abe-Ouchi, A., Braconnot, P., Brewer, S., Kageyama, M., Kitoh, A., Legrande, A., Marti, O., et al.: Mid-Holocene NAO: A PMIP2 model intercomparison, *Geophys. Res. Lett.*, 32, 2005.
- 35 Harrison, S., Bartlein, P., Brewer, S., Prentice, I., Boyd, M., Hessler, I., Holmgren, K., Izumi, K., and Willis, K.: Climate model benchmarking with glacial and mid-Holocene climates, *Clim. Dynam.*, 43, 671–688, 2014.

- Heuvelink, G. et al.: Propagation of error in spatial modelling with GIS, *Geogr. Inf. Syst.*, 1, 207–217, 1999.
- Hickler, T., Vohland, K., Feehan, J., Miller, P. A., Smith, B., Costa, L., Giesecke, T., Fronzek, S., Carter, T. R., Cramer, W., et al.: Projecting the future distribution of European potential natural vegetation zones with a generalized, tree species-based dynamic vegetation model, *Global. Ecol. Biogeogr.*, 21, 50–63, 2012.
- 5 Kaplan, J. O., Krumhardt, K. M., and Zimmermann, N.: The prehistoric and preindustrial deforestation of Europe, *Quaternary. Sci. Rev.*, 28, 3016–3034, 2009.
- Klein Goldewijk, K., Beusen, A., Van Drecht, G., and De Vos, M.: The HYDE 3.1 spatially explicit database of human-induced global land-use change over the past 12,000 years, *Global. Ecol. Biogeogr.*, 20, 73–86, 2011.
- Koenigk, T., Brodeau, L., Graversen, R. G., Karlsson, J., Svensson, G., Tjernström, M., Willén, U., and Wyser, K.: Arctic climate change in
 10 21st century CMIP5 simulations with EC-Earth, *Clim. Dynam.*, 40, 2719–2743, 2013.
- Lindgren, F., Håvard, R., and Lindström, J.: An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach, *J. Roy. Statist. Soc. Ser. B*, 73, 423–498, 2011.
- Mikolajewicz, U., Gröger, M., Maier-Reimer, E., Schurgers, G., Vizcaíno, M., and Winguth, A. M.: Long-term effects of anthropogenic CO₂ emissions simulated with a complex earth system model, *Clim. Dynam.*, 28, 599–633, 2007.
- 15 Miller, P. A. and Smith, B.: Modelling tundra vegetation response to recent arctic warming, *Ambio*, 41, 281–291, 2012.
- Olofsson, J.: The Earth: climate and anthropogenic interactions in a long time perspective, Lund University, 2013.
- Päivinen, R., Lehtikoinen, M., Schuck, A., Häme, T., Väättäinen, S., Kennedy, P., and Folving, S.: Combining earth observation data and forest statistics, *EuroForIns*, 2001.
- Pirzamanbein, B., Lindström, J., Poska, A., Sugita, S., Trondman, A.-K., Fyfe, R., Mazier, F., Nielsen, A. B., Kaplan, J. O., Bjune, A. E.,
 20 Birks, H. J. B., Giesecke, T., Kangur, M., Latałowa, M., Marquer, L., Smith, B., and Gaillard, M.-J.: Creating spatially continuous maps of past land cover from point estimates: A new statistical approach applied to pollen data, *Ecol. Complex.*, 20, 127 – 141, 2014.
- Pirzamanbein, B., Lindström, J., Poska, A., and Gaillard, M.-J.: Modelling Spatial Compositional Data: Reconstructions of past land cover and uncertainties, *arXiv preprint arXiv:1511.06417*, 2015.
- Pitman, A., de Noblet-Ducoudré, N., Cruz, F., Davin, E., Bonan, G., Brovkin, V., Claussen, M., Delire, C., Ganzeveld, L., Gayler, V., et al.:
 25 Uncertainties in climate responses to past land cover change: First results from the LUCID intercomparison study, *Geophys. Res. Lett.*, 36, 2009.
- Pongratz, J., Reick, C., Raddatz, T., and Claussen, M.: A reconstruction of global agricultural areas and land cover for the last millennium, *Global. Biogeochem. Cy.*, 22, 2008.
- Prentice, I. C., Bondeau, A., Cramer, W., Harrison, S. P., Hickler, T., Lucht, W., Sitch, S., Smith, B., and Sykes, M. T.: Dynamic global
 30 vegetation modeling: quantifying terrestrial ecosystem responses to large-scale environmental change, in: *Terrestrial ecosystems in a changing world*, pp. 175–192, Springer, 2007.
- Richter-Menge, J. A., Jeffries, M. O., and Overland, J. E.: Arctic report card 2011, National Oceanic and Atmospheric Administration, 2011.
- Rue, H. and Held, L.: Gaussian Markov random fields: theory and applications, CRC Press, 2004.
- Samuelsson, P., Jones, C. G., Willén, U., Ullerstig, A., Gollvik, S., Hansson, U., Jansson, C., Kjellström, E., Nikulin, G., and Wyser, K.: The
 35 Rossby Centre Regional Climate model RCA3: model description and performance, *Tellus. A*, 63, 4–23, 2011.
- Scheiter, S., Langan, L., and Higgins, S. I.: Next-generation dynamic global vegetation models: learning from community ecology, *New. Phytol.*, 198, 957–969, 2013.

- Schuck, A., van Brusselen, J., Päivinen, R., Häme, T., Kennedy, P., and Folving, S.: Compilation of a calibrated European forest map derived from NOAA-AVHRR data, EFI Internal Report 13, EuroForIns, 2002.
- Sitch, S., Smith, B., Prentice, I. C., Arneth, A., Bondeau, A., Cramer, W., Kaplan, J., Levis, S., Lucht, W., Sykes, M. T., et al.: Evaluation of ecosystem dynamics, plant geography and terrestrial carbon cycling in the LPJ dynamic global vegetation model, *Glob. Change Biol.*, 9, 161–185, 2003.
- Smith, B., Prentice, I. C., and Sykes, M. T.: Representation of vegetation dynamics in the modelling of terrestrial ecosystems: comparing two contrasting approaches within European climate space, *Global. Ecol. Biogeogr.*, 10, 621–637, 2001.
- Strandberg, G., Brandefelt, J., Kjellström, E., and Smith, B.: High-resolution regional simulation of last glacial maximum climate in Europe, *Tellus. A*, 63, 107–125, 2011.
- 10 Strandberg, G., Kjellström, E., Poska, A., Wagner, S., Gaillard, M.-J., Trondman, A.-K., Mauri, A., Davis, B. A. S., Kaplan, J. O., Birks, H. J. B., Bjune, A. E., Fyfe, R., Giesecke, T., Kalnina, L., Kangur, M., van der Knaap, W. O., Kokfelt, U., Kuneš, P., Latałowa, M., Marquer, L., Mazier, F., Nielsen, A. B., Smith, B., Seppä, H., and Sugita, S.: Regional climate model simulations for Europe at 6 and 0.2 k BP: sensitivity to changes in anthropogenic deforestation, *Clim. Past.*, 10, 661–680, 2014.
- 15 Trondman, A.-K., Gaillard, M.-J., Mazier, F., Sugita, S., Fyfe, R., Nielsen, A. B., Twiddle, C., Barratt, P., Birks, H. J. B., Bjune, A. E., Björkman, L., Broström, A., Caseldine, C., David, R., Dodson, J., Dörfler, W., Fischer, E., van Geel, B., Giesecke, T., Hultberg, T., Kalnina, L., Kangur, M., van der Knaap, P., Koff, T., Kuneš, P., Lagerås, P., Latałowa, M., Lechterbeck, J., Leroyer, C., Leydet, M., Lindbladh, M., Marquer, L., Mitchell, F. J. G., Odgaard, B. V., Peglar, S. M., Persson, T., Poska, A., Rösch, M., Seppä, H., Veski, S., and Wick, L.: Pollen-based quantitative reconstructions of Holocene regional vegetation cover (plant-functional types and land-cover types) in Europe suitable for climate modelling, *Glob. Change Biol.*, 21, 676–697, 2015.
- 20 Zhang, W., Miller, P. A., Smith, B., Wania, R., Koenigk, T., and Döscher, R.: Tundra shrubification and tree-line advance amplify arctic climate warming: results from an individual-based dynamic vegetation model, *Environ. Res. Lett.*, 8, 034 023, 2013.