

# Statistically-estimated tree composition for the northeastern United States at the time of European settlement

December 11, 2014

Christopher J. Paciorek, Andrew Thurman, Simon J. Goring, Charlie V. Cogbill,  
John W. Williams, David J. Mladenoff, Jody A. Peters, Jun Zhu,  
and Jason S. McLachlan [author order to be discussed]  
[please tell me if you do/don't want your middle initial]

TODO:

- look through code and labbook to make sure have all details
- properly cite Goring stuff and reference details of data collection
- focus on PLOS One, but Nature's Scientific Data is a possibility

## Abstract

We present a data product of the estimated composition of tree taxa at the time of European settlement of the United States and the statistical methodology used to produce the product. Composition is defined as the proportion of stems larger than approximately X cm diameter at breast height [Charlie? Simon? is X roughly 10?] in 22 taxonomic groupings, generally at the genus level. The data come from settlement survey records that provide raw data that are transcribed and then aggregated spatially, giving count data. The domain is divided into two regions, western (Indiana through Minnesota) and eastern (Ohio to Maine). Public Land Survey point data in the western region is aggregated to a regular 8 km grid, while data in the eastern region, from Town Proprietor Surveys is aggregated at the township level. The product is based on a Bayesian statistical model fit to the count data that estimates composition on a regular 8 km grid. The statistical model allows us to estimate composition at locations with no data and to smooth over noise caused by limited counts in locations with data. Critically, it also allows us to quantify uncertainty in our composition estimates. We expect this data product to be useful for understanding the state of vegetation in the northeastern United States prior to large-scale European settlement. In addition to specific regional questions, the data product can also serve as a baseline against which to investigate how forests and ecosystems change after intensive settlement.

# 1 Introduction

Historical datasets provide critical context to understand forest ecology. Historical datasets allow researchers to define 'baseline' conditions in conservation management, to understand ecosystem processes at decadal and centennial scales, and, particularly in regions with widespread land use change, to understand the extent to which forest conversion and regeneration differ from the original forest cover.

Euro-American settlement and subsequent land use change occurred in a time transient fashion across North America, and surveys of vegetation followed a similar pattern. Early surveys (from 1620 until 1825) in the northeastern United States were areal surveys at a township level (Cogbill et al., 2002; Thompson et al., 2013), which we call the Town Proprietor Survey (TPS). These surveys reporting percent composition of major tree species using sometimes inconsistent common names. Later surveys, after the establishment of the Public Land Survey System (PLSS) - from 1819 to 1904 - were point estimates along a regular grid, with 1 mile spacing (Schulte and Mladenoff, 2001; Bourdo, 1956; Goring et al., prep). Survey instructions during the PLSS varied through time, and by point type. This required the application of spatially varying correction factors (Cogbill et al., prep; Goring et al., prep) to accurately assess stem density, basal area and biomass from the early settlement records.

Logging, agriculture and abandonment has left an indelible mark on forests in the northeastern United States (Foster et al., 1998; Goring et al., prep; Rhemtulla et al., 2009b; Thompson et al., 2013), however most studies have studied these effects at the state level or below (Rhemtulla et al., 2009a; Friedman and Reich, 2005), and with varying spatial resolution, from townships (36 square miles) to forest zones of hundreds or thousands of square miles. Goring et al. (prep) provide a dataset across the upper Midwest that is resolved to an 8 x 8 km grid size, providing broad spatial coverage at a meso-spatial scale that can be compared to the Forest Inventory and Analysis products Gray et al. (2012). The limitation of this dataset is that it provides only estimates of variance within cell for PLSS data in the western region, and no township estimates of error for regions surveyed as part of the Town Proprietor Survey.

Properly assessing uncertainty in ecological data is imperative to understanding and modelling ecological processes (Cressie et al., 2009). In this way, a model that can account for the spatial structure of the underlying PLSS and TPS data, and provide reliable estimates of uncertainty across the northeastern United States provides an enormously valuable tool for researchers interested in the ecological structure and function of forests at longer time scales.

# 2 Data

The raw data are obtained from survey records collated from across the northeastern U.S. by a number of researchers. For the states of Minnesota, Wisconsin, Illinois, Indiana, and Michigan (the western subdomain), data are available at survey point locations and have been aggregated to a regular 8 km grid in the Albers projection. For the states of Ohio, Pennsylvania, New Jersey, New York and the six New England states (the eastern subdomain), data are aggregated at the township level. There is also data from a single township in Quebec and a single township in northern Delaware. Data are essentially complete in Minnesota, Wisconsin and Michigan but data in Illinois and Indiana represent a sample of the full set of grid cells, with survey record

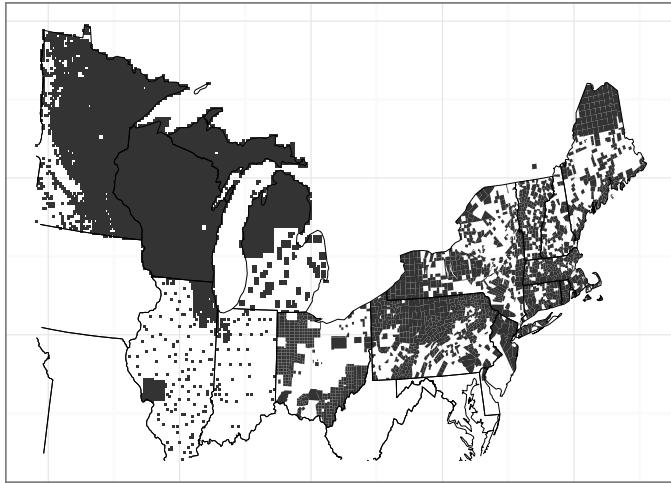


Figure 1: Spatial domain, with locations with data shown in gray. Locations are grid cells in western portion and townships in eastern portion.

transcription ongoing. Data for the remaining states are a subset of the full set of townships. Fig. 2 [Fig. 1 - Lyx is having a numbering problem] shows the domain, indicating the grid cells and townships with data.

Note that surveys occurred over a period of more than 200 years as European colonists (before U.S. independence) and the United States settled what is now the northeastern United States. Our estimates are for the period of settlement represented by the survey data and therefore are time-transgressive; they do not represent any single point in time across the domain, but rather the state of the landscape at the time just before settlement occurred [get ref from Charlie or Jason's historian pal]

Extensive details on the data are available in Goring et al. (2015) and the raw data are available at XXX. The aggregation into taxonomic groups is primarily at the genus level, but in some cases of monospecific genera, is at the species level. We model the following 22 taxa plus an “other hardwood” category: Atlantic white cedar (*Chamaecyparis thyoides*), Ash (*Fraxinus spp.*), Basswood (*Tilia americana*), Beech (*Fagus grandifolia*), Birch (*Betula spp.*), Black gum/sweet gum (*Nyssa sylvatica* and *Liquidambar styraciflua*), Cedar/juniper (what is in here other than *Juniperus* and is the *Juniperus* only *Juniperus virginiana*), Cherry (*Prunus spp.*), Chestnut (*Castanea dentata*), Dogwood (*Cornus spp.*), Elm (*Ulmus spp.* or *Ulmus americana???*), Fir (*Abies spp.*), Hemlock (*Tsuga canadensis*), Hickory (*Carya spp.*), Ironwood (*Carpinus caroliniana* and *Ostrya virginiana* (anything else?)), Maple (*Acer spp.*), Oak (*Quercus spp.*), Pine (*Pinus spp.*), Poplar/tulip poplar (*Populus spp.* and *Liriodendron tulipifera*), Spruce (*Picea spp.*), Tamarack (*Larix laricina*), Walnut (*Juglans nigra*). [Simon/Charlie to check my Latin and whether I've left out anything that falls into a given common name] Note that in several cases (e.g., black gum/sweet gum, ironwood, poplar/tulip poplar, cedar/juniper), because of ambiguity in the common tree names used by surveyors, a group represents trees from different families and even orders. For the western subdomain we do not fit statistical models for Atlantic white cedar and chestnut as these have 0

and 5 trees present in the dataset, respectively.

[include description of size cutoff] [Simon, Jody mentioned that ther are 402 trees in IN below ~10 cm and at least 43 in IL. Do you screen those out or do we just live with it?]

There are approximately 520,000 trees from the western subdomain and 420,000 trees from the eastern subdomain. In each subdomain, oak is the most common taxon and pine the second most common.

Our domain is a rectangle covering all of the states, in the Albers projection (NAD 1983 Great Lakes and St. Lawrence Albers) [Simon, please confirm this language is right], with the rectangle split into 8 km cells, arranged in a 296 by 180 grid of cells, with the centroid of the cell in the southeast corner located at (-71000, 58000). For the modeling of the western subdomain we use the western-most 146 by 180 grid of cells. For the modeling of the eastern subdomain we use the eastern-most 180 by 180 grid of cells and then omit 23 cells in the north and 17 cells in the south. [check 23 in N vs 17 in N; also check what is in netCDF product - should be for the full domain with missing values?]

## 3 Statistical model

We fit a Bayesian statistical model to the data, with two primary goals:

1. To estimate composition on a regular grid across the entire domain, filling gaps where no data are available, and
2. To quantify uncertainty in composition at all locations. Even in grid cells and townships with data, we wish to quantify uncertainty because the empirical proportions represent estimates of the true proportions based on the full population of all the trees in an areal region.

The result of fitting the Bayesian model via Markov chain Monte Carlo (MCMC) is a set of representative samples from the posterior distribution for the composition of the taxa at all of the grid cells. These samples are the data product and can then be used in analyses. We also provide the mean and standard deviation of the samples, which represent our best estimate of composition and a Bayesian “standard error” for the estimate.

### 3.1 Data model

The statistical model treats the observations as having a multinomial distribution with a (latent) vector of proportions for each grid cell.

$$y_i \sim \text{Multi}(n_i, \theta(s_i))$$

where  $y_i$  is the vector of counts for the  $P$  taxa at the  $i$ th location,  $n_i$  is the number of trees counted in the location, and  $\theta(s_i)$  is the vector of unknown proportions for those taxa at that location. Note that we use a standard multinomial distribution without overdispersion as the set of trees in the dataset is roughly uniformly sampled across the cells or townships (Goring et al., 2015).

In what follows we describe the basic model for those states for which we have raw data on the 8 km grid and in Section 3.5 we describe the extension of the model to accommodate data aggregated to the township level.

The proportions,  $\theta_p(s_i)$ ,  $p = 1, \dots, P$ , are modeled spatially by a set of  $P$  Gaussian spatial processes, one per taxon,  $\alpha_p(s_i)$ ,  $p = 1, \dots, P$ . This collection of processes defines a multivariate spatial process for composition. The  $\alpha_p(s)$  processes are defined on the 8 km grid,  $\alpha_p = \{\alpha_p(s_1), \dots, \alpha_p(s_m)\}$  for the  $m$  grid cells. In Section 3.4 we introduce a multinomial probit model that relates the  $\alpha_p(s)$  processes to the proportion processes,  $\theta_p(s)$ , via the introduction of latent variables, with an implicit sum-to-one constraint,  $\sum_{p=1}^P \theta_p(s) = 1$ .

The critical component of the statistical model is the representation of  $\alpha_p(s)$  as a spatial process. This process is a prior structure that serves to smooth across noise in the observations and allows for interpolation to locations with no data. Apart from the sum to one constraints, the taxa are considered to be independent in the prior, as we did not want to impose any structure that ties the different taxa together, as any correlation will likely vary across space.

In the next section, we consider two spatial models to define the structure of the  $\alpha_p(s)$  processes, a standard conditional autoregressive model (Banerjee et al., 2003) and a Gaussian Markov random field (MRF) approximation to a Gaussian process with Matern covariance (Lindgren et al., 2011).

## 3.2 Spatial process models

MRF models work directly with the precision matrix of the values of the spatial process, so calculation of the prior density of  $\alpha_p(s)$  is computationally simple (Rue and Held, 2005), but in situations where the likelihood is not normal, it can be difficult to set up effective MCMC algorithms that are able to move in the high-dimensional space of  $\alpha_p$ . The latent variable representation of Section 3.5 helps to alleviate this problem.

### 3.2.1 Standard conditional autoregressive models

Our first model is a standard conditional autoregressive (CAR) model (Banerjee et al., 2003). We use a standard form of this model, which treats the four cardinal neighbors of each grid cell as the neighbors of the grid cell. The corresponding precision matrix has diagonal elements,  $Q_{ii}$ , equal to the number of neighbors for the  $i$ th area (i.e., four except for cells on the boundary of the domain), while  $Q_{ik} = -1$  (the negative of a weight of one) when areas  $i$  and  $k$  are neighbors and  $Q_{ik} = 0$  when they are not. This gives the following model for the values of  $\alpha_p(s_i)$  collected as a vector across all of the grid cells,  $i = 1, \dots, m$ :

$$\alpha_p \sim N(0, \sigma_p^2 Q^-)$$

The use of the generalized inverse notation indicates that  $Q$  is not full-rank, but is of rank  $m - 1$ ; this gives an improper prior on an implicit overall mean for the process values. This specification is called an *intrinsic conditional autoregression (ICAR)* and  $Q$  can also be expressed as  $D - C$  where  $C$  is the  $m \times m$  adjacency matrix defining the neighborhood relation of the locations; that is  $(C)_{ik} = 1$  if locations  $i$  and  $k$  are neighbors. The matrix  $D$  is an  $m \times m$  diagonal matrix containing the row sums of matrix  $C$  as the diagonal entries  $(D)_{ii} = \sum_{k=1}^m (C)_{ik}$ .

We refer to this as the *CAR model*.

### 3.2.2 Gaussian process approximation

Gaussian processes (GP) are also standard models for spatial processes (Banerjee et al., 2003). GP models are computationally challenging for large datasets because of manipulations involving large covariance matrices. Given this, Lindgren et al. (2011) proposed a new framework for using Gaussian MRFs (GMRFs) as approximations to GPs, based on the use of stochastic partial differential equations (SPDE).

We consider Gaussian processes in the Matérn class, using the following parameterization of the Matérn correlation function as

$$R(d) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left( \frac{2\sqrt{\nu}d}{\rho} \right)^{\nu} \mathcal{K}_{\nu} \left( \frac{2\sqrt{\nu}d}{\rho} \right), \quad (1)$$

where  $d$  is Euclidean distance,  $\rho$  is the spatial range parameter, and  $\mathcal{K}_{\nu}(\cdot)$  is the modified Bessel function of the second kind, whose order is the smoothness (differentiability) parameter,  $\nu > 0$ .  $\nu = 0.5$  gives the exponential covariance. For any pair of locations,  $R(d)$  defines the correlation of the process, e.g.,  $\alpha_p(s)$ , as a function of the distance between the locations. Considering all pairs of locations, this defines a correlation matrix for all locations of interest.

The Lindgren et al. (2011) approach allows us to consider MRF approximations to the Matern-based GPs for  $\nu = 1$  and  $\nu = 2$ . Our second model is the Lindgren approximation for Matern-based GPs with  $\nu = 1$ . To implement the Lindgren model, one modifies the  $Q$  matrix defined previously as follows. Let  $a = 4 + \frac{1}{\rho^2}$ . The diagonal elements of  $Q$  are  $4 + a^2$ . The entries corresponding to cardinal neighbors are  $-2a$ . Those for diagonal neighbors are 2 and those for 2nd-order cardinal neighbors are 1. This extends the neighborhood structure relative to the CAR model and parameterizes it as a function of  $\rho$ .

The primary difference between the CAR and Lindgren models is that the Lindgren model provides an additional degree of freedom by estimating  $\rho$ . In particular  $\rho$  allows us to estimate the locality of the smoothing. As  $\rho$  decreases, the model uses more and more local data to estimate the compositional proportions at a given location, effectively averaging the empirical proportions over smaller neighborhoods. In general, the Lindgren et al. (2011) model will generally provide for a smoother estimate than the CAR model (Paciorek, 2013).

To ensure that the  $\sigma^2$  parameter is equivalent between the two models, we reparameterize, producing our second model:

$$\alpha_p \sim N(\mu_p, \sigma_p^2 \cdot \frac{4\pi}{\rho_p^2} Q(\rho_p)^{-1})$$

We refer to this model as the *SPDE model*.

### 3.3 Prior Distributions

The ICAR specification contains a set of hyperparameters  $\sigma_p^2$  for  $p = 1, \dots, P$ . Following (Gelman, 2006) we use a uniform distribution on each  $\sigma_p$  parameter, with upper bound of [check max values]. For the SPDE model we also have parameters  $\mu_p$ , which we give flat, non-informative priors (truncated at  $\pm 10$ ), and  $\rho_p$  which we give uniform priors on the interval  $(0.1, \exp(5))$ .

### 3.4 Latent Variable Model

It is well-known that devising an effective MCMC algorithm for models with latent Gaussian process(es) and a non-Gaussian likelihood is difficult (cite recent cent/noncent papers). To develop an algorithm, we make use of a latent variable representation for the multinomial probit model (McCulloch and Rossi, 1994). The representation introduces latent variables that allow one to develop a MCMC sampling strategy that takes advantage of closed form full conditional distributions (so-called Gibbs sampling steps) for  $\alpha_p$ .

Suppose that compositional counts are available at a number of locations. At location  $i$ , a sample size of  $n_i$  observations are collected, and each observation (i.e., each tree) can be classified into  $P$  distinct categories. For a given tree  $j$  at location  $i$ , let  $Y_{ij}$  denote the response variable indicating the category. Let  $Y_{ij}$  be associated with  $P$  latent variables  $W_{ij1}, \dots, W_{ijP}$  such that  $Y_{ij} = p$  if and only if  $W_{ijp} = \max_{p'} \{W_{ijp'}\}$ ; in other words, the maximum of the set of latent variables  $\{W_{ijp}\}_{p=1}^P$  determines the category of observation  $j$  at location  $i$ .

The final piece of the latent variable representation is the relationship between the  $W$  variables and the  $\alpha_p(s)$  processes. We have that

$$W_{ijp} \sim N(\alpha_p(s_i), 1)$$

independently for all of the  $W_{ijp}$  values.

Consider the following example with two locations that are neighbors and  $P = 2$  categories. Each tree  $j$  at location  $i$  is associated with two variables  $W_{ij1}$  and  $W_{ij2}$ , governed by the latent variables  $\alpha_{i1}$  and  $\alpha_{i2}$ , respectively. Suppose that  $\alpha_{i1} > \alpha_{i2}$  for a given location  $i$ . Then this model implies that any tree  $j$  is more likely to be labeled 1 than 2 at location  $i$ . The difference between  $\alpha_{i1}$  and  $\alpha_{i2}$  explains the *difference* in probability of *categories* 1 and 2 at location  $i$ , and the similarity between  $\alpha_{1p}$  and  $\alpha_{2p}$  explains the *correlation* between the probabilities at *locations* 1 and 2 for category  $p$ .

### 3.5 Model for township data

We developed an extension of the model described in previous sections to account for data at a different aggregation than our core 8 km grid. This extension introduces a new set of latent variables, one per tree, that indicate the grid cell in which the tree is located. These latent 'membership' variables,  $c_{tj}$ , for tree  $j$  in township  $t$  can be sampled within the MCMC as additional unknown parameters. The prior for  $c_{tj}$  is a discrete distribution that puts mass proportional to the overlap between the township in which the tree is located and the grid cells that intersect the township,  $\psi_{1(t)}, \dots, \psi_{T(t)}$ , with the priors across different trees in a township being independent.

$$c_{tj} \sim \text{Multinom}(1, \{\psi_{1(t)}, \dots, \psi_{T(t)}\}).$$

In updating the other parameters in the model, we condition on the current values,  $c_{tj}$ , which provides a "soft" assignment of trees to grid cells that respects both the township in which the tree was sampled and the uncertainty in which grid cell the tree was sampled.

Note that this prior has the unrealistic feature that it does not represent our knowledge that the trees in a township would be distributed more regularly across the area of the township than expected by such an independence prior.

### 3.6 Computation

The McCulloch and Rossi (1994) representation is convenient for MCMC sampling, particularly in this high-dimensional spatial context, as it allows us to draw from the posterior conditional distribution of the  $W_{ijp}$  variables (these distributions are truncated normal) in closed form and to draw the entire vector of latent process values,  $\alpha_p$ , as a single sample that respects the spatial dependence structure for each taxon.

While the latent variable representation provides great advantages in the MCMC sampling for each  $\alpha_p$  compared to joint Metropolis updates or single location at a time updates, there is still strong dependence between the hyperparameters,  $\{\sigma_p^2, \mu_p, \rho_p\}$  and the the latent process values (as well as between the latent process values and the latent  $W_{ijp}$  variables). To address the first, we developed a 'cross-level' joint updating strategy in which we propose  $\phi \in \{\sigma_p\}, p = 1, \dots, P$  and for the SPDE model,  $\phi \in \{\{\sigma_p\}, \{\mu_p\}, \{\rho_p\}\}$ , via a Metropolis-style random walk and then conditional on the proposed value of  $\phi$  propose  $\alpha_p$  from its full conditional distribution  $\backslash \alpha_p$  and sampling from the marginalized distribution  $\phi | W$ . Note that in sampling  $\sigma_p$  and  $\rho_p$  in the SPDE model, for each  $p$ , we jointly propose  $\{\sigma_p, \rho_p\}$  using an adaptive Metropolis proposal and then use the cross-level strategy to also propose  $\alpha_p$  as part of a single accept-reject step. For these various joint samples of hyperparameters and  $\alpha_p$ , we use adaptive Metropolis sampling (Shaby and Wells, 2011).

The full description of the MCMC sampling steps is provided in the Appendix. In addition, in the latent variable representation,  $\theta_p(s)$  never appears explicitly and cannot be calculated in close form. Instead we use Monte Carlo integration over  $W_{ijp}$ ,  $p = 1, \dots, P$  to estimate  $\theta_p(s_i)$ , also described in the Appendix.

The model is implemented in R with core computational calculations coded in C++ using the Rcpp package (Eddelbuettel and Francois, 2011). We also make extensive use of sparse matrix representations and algorithms, using the spam package in R (Furrer and Sain, 2010). All code is available on Github, including pre- and post-processing code.

## 4 Model comparison

### 4.1 Design

We compared the CAR and SPDE models using cross-validation by holding out data from the fitting process and assessing the fit of the model on the held-out data. We used two types of test data:

1. We held out all the data from 95% of the cells in Minnesota, with cells selected at random. This was meant to assess the ability of the model to interpolate from a sparse set of cells/townships and mimics the limited data in Illinois and Indiana.
2. We held out 5% of the trees from all of the trees in the dataset (leaving aside the held-out Minnesota cells). This was meant to assess the ability of the model to estimate the composition in cells in which data were available.

Finally, in a separate sensitivity analysis we instead left out 80% of the cells in Minnesota at random. This variation on the first test above was meant to indicate whether our model comparison

conclusions would be robust as the digitization process for Illinois and Indiana progresses and provides us with increasingly dense data.

There has been extensive work in the statistical literature on good metrics to use to compare the predictive ability of models; these metrics are referred to as scoring rules. In particular it is well-recognized that predictive distributions should maximize sharpness subject to calibration. That is the predictive distribution should be as narrow as possible while being calibrated such that the observations are consistent with the distribution (Gneiting et al., 2007). When thinking in terms of prediction intervals, we seek intervals that are as narrow as possible while still covering the truth the expected proportion of the time. Following the suggestions in (Gneiting et al., 2007), we consider these metrics.

We used the following criteria to compare models. For experiment 1,  $Y_i = \{Y_{i1}, \dots, Y_{iP}\}$  is the count of all trees in the held-out cell and for experiment 2,  $Y_i$  is the count of held-out individual trees in the location.

1. Brier score: (Gneiting et al., 2007) suggest this measure, which has been in use for decades. For multi-category as opposed to binary outcomes, this takes the form

$$\frac{1}{n_i} \sum_{j=1}^{n_i} \sum_{p=1}^P (y_{ijp} - \theta_p(s_i))^2$$

for each cell. The full Brier score for the complete set of held-out locations sums over all locations.

2. Log predictive density: This measure takes the log of the density of held-out observations under the fitted model,  $Y_i \sim \text{Multinom}(n_i, \{\theta_1(s_i), \dots, \theta_P(s_i)\})$ , summing on the log scale across all of the held-out data.

While in principle, this metric should be optimal (Krnjajić and Draper, 2014), it suffers from a lack of robustness in being very sensitive to small predictions near zero (Gneiting et al., 2007). Even worse, our Monte Carlo estimation of  $\theta$  used 10000 samples, so in some cases  $\theta_p(s) = 0$ . When a tree is present in a cell but its corresponding proportion is 0 this gives a log density of  $-\infty$ , preventing use of the metric. As an informal solution to this we set  $\theta_p(s) = 0.5 \frac{1}{10000}$  in such cases, but given these issues we treat the log predictive density as a secondary measure.

3. Weighted root mean square prediction error (RMSPE) and mean absolute error (MAE) for individual cells (only for experiment 1), where we weight by the number of held-out trees to account for the greater variability in the empirical proportions in locations with few held-out trees.
4. Coverage and length of 95% prediction intervals for  $Y_{ip}$ .

For each of these metrics, we calculated the metric using the posterior mean composition estimates (as a measure of our core predictions) and averaging the metric over the posterior samples (as a measure of our full data product, including uncertainty). In other words in the first we first average over the composition posterior samples and in the second we average of the posterior values of the cross-validation metric.

Table 1: Predictive ability based on several predictive score criteria for the CAR and SPDE spatial models when holding out 95% of entire cells of data in Minnesota. Smaller values are better.

	Posterior mean of score			Score of posterior mean predictions	
	CAR model	SPDE model	Posterior Prob. CAR < SPDE	CAR model	SPDE model
Brier	84997	86803	0.92	76068	75330
Negative Log Density	222144	242676	1.00	185748	185843
Mean Absolute Error	0.0382	0.0390	0.92	0.0303	0.0293
Root Mean Square Error	0.0944	0.0986	0.77	0.0701	0.0677

In our fitting we noticed that the SPDE model produced boundary effects in the predicted composition near the edges of the convex hull of the observations. To attempt to alleviate this, we added a buffer zone of 6 grid cells around our entire original domain, though the boundary effects were still evident even after inclusion of the buffer. For the model comparison, for comparability, we included this buffer for both the SPDE and CAR models.

For the model comparison, we used an earlier version of the dataset than that used for the final data product as the model comparison runs involved substantial computation. This earlier version had incomplete data for southern Michigan and was missing a small number of trees (39 individual trees) classified as “other hardwood” due to an update to our taxonomic aggregations of common names used by surveyors. Given that part of the goal of the exercise is to understand the ability of the models to interpolate to areas without data and the ongoing nature of digitization, we believe this difference in datasets is not an issue.

We ran each model for XXX iterations, retaining YYY after burnin and subsampling to reduce storage needs.

## 4.2 Results

Here we summarize the results of our cross-validation analyses that inform the choice between the CAR and SPDE models.

### 4.2.1 Full cell hold-out experiment

Table 1 shows predictive ability for Experiment 1: those cells in Minnesota held out of the fitting process, for both the posterior mean predictions and the full posterior sample. Table 2 shows performance of the uncertainty estimates. For the posterior distribution over the predictive score values, the CAR model generally outperforms the SPDE model, while for the posterior mean predictions, the SPDE model appears to outperform the CAR model, but we do not have any uncertainty estimates for this comparison. Coverage and interval lengths are similar between the two models. From a practical perspective, based on the difference in mean absolute error, which is on a readily interpretable scale, the differences between the models are small.

Table 3 shows the results when the proportion of cells that are held out decreases from 95% to 80%, indicating performance on less sparse data. Table 4 shows performance of the uncertainty estimates. For denser data, unlike the results when data are more sparse, the SPDE model generally

Table 2: Coverage and length of prediction intervals for the CAR and SPDE spatial models when holding out 95% of entire cells of data in Minnesota.

	CAR model	SPDE model
Coverage	0.977	0.978
Mean Interval Length	0.135	0.142
Median Interval Length	0.040	0.031

Table 3: Predictive ability based on several predictive score criteria for the CAR and SPDE spatial models when holding out 80% of entire cells of data in Minnesota. Smaller values are better, except for coverage, where values near 0.95 are optimal.

	Posterior mean of score			Score of posterior mean predictions	
	CAR model	SPDE model	Posterior Prob. CAR < SPDE	CAR model	SPDE model
Brier	65870	65004	0.08	60169	60020
Negative Log Density	167947	167536	0.43	145404	145750
Mean Absolute Error	0.0322	0.0307	0.08	0.0246	0.0240
Root Mean Square Error	0.0800	0.0771	0.01	0.0578	0.0571

outperforms the CAR model, but again differences from a practical perspective, based on mean absolute error, are limited.

#### 4.2.2 Individual tree hold-out experiment

Table 5 shows the predictive ability for Experiment 2. Here we have limited evidence (posterior probability of 0.84) that the SPDE model is better based on the Brier score.

#### 4.2.3 Choice of spatial model

The differences between models are not consistent across the various comparisons, so there is not a clear choice. In our final data product we use the CAR model, for three reasons. First, the CAR model has modestly better performance when data are sparse, as is still the case for Illinois and Indiana. Second, the model is simpler and easier to explain and computations can be done more quickly. Third, predictions from the SPDE model showed boundary effects, with taxa showing non-negligible posterior mean values at the edges of the domain, well away from where the taxon was present in the empirical data. This included non-negligible values within (but near the edge

Table 4: Coverage and length of prediction intervals for the CAR and SPDE spatial models when holding out 80% of entire cells of data in Minnesota.

	CAR model	SPDE model
Coverage	0.982	0.974
Mean Interval Length	0.117	0.107
Median Interval Length	0.030	0.021

Table 5: Predictive ability based on several predictive score criteria for the CAR and SPDE spatial models when holding out 5% of trees. Smaller values are better.

	Posterior mean of score			Score of posterior mean predictions	
	CAR model	SPDE model	Posterior Prob. CAR < SPDE	CAR model	SPDE model
Brier	13158	13143	0.16	13040	13039
Negative Log Density	31984	31870	0.02	31212	31229

of) the convex hull of locations with data.

## 5 Data product

For the final data product, we ran the model using the CAR specification for X iterations, discarding the first Z iterations for burn-in and retaining Y to limit storage needs and limit the size of the data product. Mixing was generally reasonable, but for some of the hyperparameters was relatively slow, particularly for less common taxa. Despite this, mixing for the variables of substantive interest, the proportions was good. [Report ESS values]

In Figure 5, we show maps of estimated composition for the full domain for several taxa of substantive interest to illustrate the results: beech, cherry, chestnut, elm, hemlock, oak, and pine. These maps contrast the raw data proportions, the posterior means and posterior standard deviations as pointwise estimates of uncertainty. Fig. 5 shows posterior means for all 22 (???) taxa.

While our data product is given for the full rectangular domain, including hypothetical vegetation over water, uncertainty is large outside of the core states with data and any analysis and inference should only be done within those core states.

The data product is publically available at XX under YY license as version 1.0 as of January 2015. The product is in the form of a netCDF-4 file, with dimensions x, y, and MCMC iteration. There is one variable per taxon. The PaleON project will continue to maintain this product, releasing new versions as additional data in Illinois, Indiana and Ohio are digitized.

## 6 Discussion

Note that digitization of data from Illinois and Indiana is ongoing and digitization of additional data from Ohio is planned as well. As a result, at some point we expect to have complete data for the western half of the domain. There will remain some missing townships in the eastern half of the domain.

Given the density of data and the limited differences seen between the ICAR and SPDE models, we expect the data product to be reasonably robust to the choice of spatial model, particularly in those areas with complete data. However, additional investigation of other statistical representations is of interest, in particular nonstationary spatial models and use of covariates. The biggest shortcoming of the current model is its inability to account for local features such as rivers (note the Minnesota River floodplain in evidence in the raw data). The current model, by using a simple

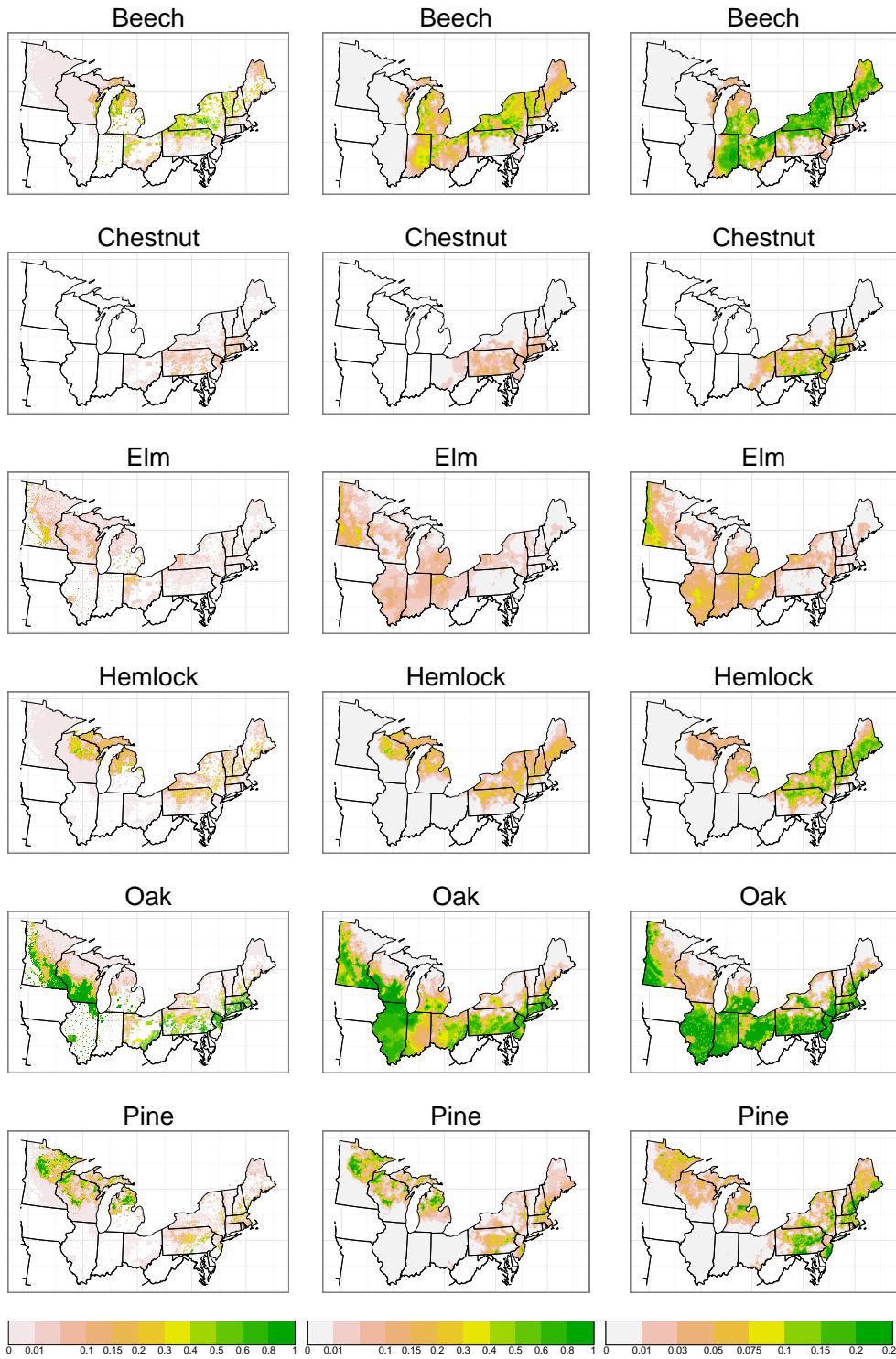


Figure 2: [THIS IS AN OLD VERSION FOR ILLUSTRATION ONLY - TO BE REPLACED WHEN VERSION 0.3 FITS ARE DONE] Empirical proportions from raw data (column 1), posterior mean predictions (column 2) and posterior standard deviations – representing standard errors of prediction (column 3) for select taxa.



Figure 3: [THIS IS AN OLD VERSION FOR ILLUSTRATION ONLY - TO BE REPLACED AND PUT IN 8 ROW X 3 COLUMN FORM WHEN VERSION 0.3 FITS ARE DONE] Posterior mean predictions for all taxa over the entire domain.

stationary spatial model that smooths as a function of Euclidean distance, does not account for topographic, soil, or other features.

Other related data products that are under development include:

- The gridded raw count data for the western subdomain and data aggregated to township for the eastern subdomain. The raw data product is available at XXX and that product provides the input data for the product described in this work. Note that we cannot provide the raw data at point locations because of limitations specified in the data use agreements under which we have access to the data.
- Gridded raw biomass estimates for Minnesota, Wisconsin, and Michigan [right Simon, or also for IL and IN] based on the PLS data.
- Statistically estimated biomass for Minnesota, Wisconsin, and Michigan [also IL and IN?] using a statistical model applied to the raw biomass estimates.

An additional drawback of the product is its focus on composition, which does not directly tell us about vegetation structure, in particular does not distinguish between closed forest, savanna, and prairie, of particular note in Minnesota, Wisconsin, Illinois, and into Indiana. The statistically-estimated biomass product mentioned just above will directly inform questions about vegetation structure. Extensions of that product will also estimate basal area and stem density.

## Acknowledgments

The authors are deeply indebted to all of the researchers over the years who have preserved, collected, and digitized survey records, in particular Ed Schools (Michigan DNR), { Brugam, ??? - Jody/Charlie/David/Jason/Simon please add names here}. This work was carried out by the PALEON Project with support from the National Science Foundation MacroSystems Program through grants EF-1065656, EF-1241868, DEB-1241874 and DEB-1241868.

## References

- Banerjee, S., A. Gelfand, and C. Sirmans (2003). Directional rates of change under spatial process models. *Journal of the American Statistical Association* 98(464), 946–954.
- Bourdo, E. A. (1956). A review of the general land office survey and of its use in quantitative studies of former forests. *Ecology*, 754–768.
- Cogbill, C. V., J. Burk, and G. Motzkin (2002). The forests of presettlement new England, USA: spatial and compositional patterns based on town proprietor surveys. *Journal of Biogeography* 29(10-11), 1279–1304.
- Cogbill, C. V., S. J. Goring, and A. Thurman (in prep). Estimation of robust correction factors for public land survey data.

- Cressie, N., C. A. Calder, J. S. Clark, J. M. V. Hoef, and C. K. Wikle (2009). Accounting for uncertainty in ecological analysis: the strengths and limitations of hierarchical statistical modeling. *Ecological Applications* 19(3), 553–570.
- Eddelbuettel, D. and R. Francois (2011). Rcpp: Seamless R and C++ integration. *Journal of Statistical Software* 40, 1–18.
- Foster, D. R., G. Motzkin, and B. Slater (1998). Land-use history as long-term broad-scale disturbance: regional forest dynamics in central new England. *Ecosystems* 1(1), 96–119.
- Friedman, S. K. and P. B. Reich (2005). Regional legacies of logging: departure from presettlement forest conditions in northern minnesota. *Ecological applications* 15(2), 726–744.
- Furrer, R. and S. R. Sain (2010). spam: A sparse matrix R package with emphasis on MCMC methods for Gaussian Markov random fields. *Journal of Statistical Software* 36(10), 1–25.
- Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian Analysis* 1(3), 515–534.
- Gneiting, T., F. Balabdaoui, and A. E. Raftery (2007). Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69(2), 243–268.
- Goring, S., J. Wlliams, D. Mladenoff, C. Cogbill, S. Record, C. Paciorek, S. Jackson, M. Dietze, J. Matthes, and J. McLachlan (2015). Changes in forest composition, stem density, and biomass from the settlement era to present in for the upper midwestern united states. *in preparation* 0, 0.
- Goring, S. J., J. W. Williams, D. J. Mladenoff, C. V. Cogbill, S. Record, C. J. Paciorek, S. T. Jackson, M. C. Dietze, J. H. Matthes, and J. S. McLachlan (in prep). Changes in forest composition, stem density, and biomass from the settlement era to present in for the upper midwestern united states. *Ecological Monographs*.
- Gray, A. N., T. J. Brandeis, J. D. Shaw, W. H. McWilliams, P. D. Miles, et al. (2012). Forest inventory and analysis database of the United States of America (fia). *Vegetation databases for the 21st century.–Biodiversity & Ecology* 4, 255–264.
- Krnjajić, M. and D. Draper (2014). Bayesian model comparison: Log scores and DIC. *Statistics and Probability Letters* 88, 9–14.
- Lindgren, F., H. Rue, and J. Lindström (2011). An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach. *Journal of the Royal Statistical Society B* 73, 423–498.
- McCulloch, R. and P. E. Rossi (1994). An exact likelihood analysis of the multinomial probit model. *Journal of Econometrics* 64(1), 207–240.
- Paciorek, C. (2013). Spatial models for point and areal data using markov random fields on a fine grid. *Electronic Journal of Statistics* 7, 946–972.

- Rhemtulla, J. M., D. J. Mladenoff, and M. K. Clayton (2009a). Historical forest baselines reveal potential for continued carbon sequestration. *Proceedings of the National Academy of Sciences* 106(15), 6082–6087.
- Rhemtulla, J. M., D. J. Mladenoff, and M. K. Clayton (2009b). Legacies of historical land use on regional forest composition and structure in wisconsin, USA (mid-1800s-1930s-2000s). *Eco-logical Applications* 19(4), 1061–1078.
- Rue, H. and L. Held (2005). *Gaussian Markov Random Fields: Theory and Applications*. Boca Raton: Chapman & Hall.
- Schulte, L. A. and D. J. Mladenoff (2001). The original us public land survey records: their use and limitations in reconstructing presettlement vegetation. *Journal of Forestry* 99(10), 5–10.
- Shaby, B. and M. Wells (2011). Exploring an adaptive Metropolis algorithm. Technical Report 2011-14, Department of Statistics, Duke University.
- Thompson, J. R., D. N. Carpenter, C. V. Cogbill, and D. R. Foster (2013). Four centuries of change in northeastern United States forests. *PloS one* 8(9), e72540.

## 7 Appendix

### 7.1 MCMC details

Define  $\bar{w}_{i,p} = n_i^{-1} \sum_{j=1}^{n_i} W_{ijp}$  be the average of the  $W$  values for the  $p$ th taxon in the  $i$ th grid cell.

Let  $A$  be a diagonal matrix where  $A_{ii}$  is the number of trees in the  $i$ th grid cell. When there are no trees in a grid cell,  $\bar{w}_{i,p} = 0$  and  $A_{ii} = 0$ . For the township data, at each iteration, based on the current values of the grid cell membership variables,  $c_{tj}$ , trees are aggregated into grid cells and the calculations above can then be carried out.

The conditional distribution for  $W_{ijp}$  given the other unknowns in the model and the data is as follows. Let  $\text{TN}(a, b, \mu, \tau^2)$  denote the truncated normal distribution with mean parameter  $\mu$  and variance parameter  $\tau^2$ , truncated below by  $a$  and above by  $b$ .

$$W_{ijp} \sim \begin{cases} \text{TN}\left(\max_{p^* \neq y_{ij}} w_{ijp^*}, \infty, \alpha_{y_{ij}}(s_i), 1\right), & \text{if } p = y_{ij} \\ \text{TN}\left(-\infty, w_{ijy_{ij}}, \alpha_p(s_i), 1\right), & \text{if } p \neq y_{ij} \end{cases} \quad (2)$$

In essence, the truncation value is determined by the taxon of the  $j$ th tree. For a given  $p$ , the  $W$  values for all trees in all cells can be sampled in parallel.

The conditional distribution of  $\alpha_p$  is

$$\alpha_p \sim \mathcal{N}\left(\left(A + Q_p\right)^{-1} A \bar{w}_p, \left(A + Q_p\right)^{-1}\right). \quad (3)$$

where  $Q_p = (\sigma_p^2)^{-1}Q$  for the CAR model and  $\left(\sigma_p^2 \cdot \frac{4\pi}{\rho_p^2}\right)^{-1}Q(\rho_p)$  for the SPDE model. For each hyperparameter,  $\phi \in \{\{\mu_p\}, \{\log \sigma_p, \rho_p\}\}$ , we sampled  $\{\phi, \alpha_p\}$  jointly, proposing  $\phi$  as a random walk and, conditional on the proposed value of  $\phi$ , sampling  $\alpha_p$  from the distribution just above. The joint proposal is accepted or rejected as a standard Metropolis-Hastings proposal, with adaptation of the proposal (co)variance (Shaby and Wells, 2011). As mentioned previously, for the SPDE model, we jointly proposed  $\phi = \{\log \sigma_p, \rho_p\}$  from a bivariate normal distribution to account for the posterior dependence of these parameters, while for the CAR model,  $\phi \in \{\{\log \sigma_p\}\}$ .

For the township-level data, for a given tree, we draw the latent tree membership variable,  $c_{tj}$ , from a discrete distribution by normalizing the weights,  $\{\psi_{1(t)}L_{1(t)}, \dots, \psi_{T(t)}L_{T(t)}\}$  where  $L_k$  is density of  $W_{tj}$ . if  $c_{tj} = k$ , namely the product of independent normal densities,  $W_{tjp} \sim N(\alpha_p(s_k), 1)$ , over  $p = 1, \dots, P$ . Thus the posterior weights the prior based on how consistent the  $W$  values are with the  $\alpha$  values for the candidate grid cells.

## 7.2 Estimating $\theta_p(s)$ via Monte Carlo integration

In the latent variable representation,  $\theta_p(s)$  never appears explicitly and cannot be calculated in close form. Instead we use Monte Carlo integration over  $W_{ijp}$ ,  $p = 1, \dots, P$  to estimate  $\theta_p(s_i)$ . The quantity  $\theta_p(s_i) = P(W_{ijp} = \max_{p^*} W_{ijp^*})$  defines the probability of taxon  $p$  at grid cell  $i$ . This requires one to choose the number of Monte Carlo samples, which we set at 10000. For each of the saved MCMC samples,  $k = 1, \dots, K$ , we estimate  $\hat{\theta}_p^{(k)}(s_i)$  numerically. Specifically, for  $t = 1, \dots, 10000$ , we independently draw

$$W_{ijpt}^{(k)} \sim N(\alpha_p^{(k)}(s_i), 1), p = 1, \dots, P$$

and estimate

$$\hat{\theta}_p^{(k)}(s_i) = \frac{1}{10000} \sum_{t=1}^{10000} 1(W_{ijpt}^{(k)} = \max_{p^*} W_{ijp^*t}^{(k)})$$

where  $1(\cdot)$  is the indicator function that evaluates to 1 if the expression is true and 0 if false. In other words, we calculate the proportion of times that the maximum of  $W_{ijp}$ ,  $p = 1, \dots, P$  corresponds to taxon  $p$ . Considering  $\hat{\theta}_p^{(k)}(s_i)$ ,  $k = 1, \dots, K$ , we have a sample from the posterior of  $\theta_p(s_i)$ .