

Modeling PaLEON biomass

Wesley Brooks

UW-Madison

May 28, 2013

Outline

Goal

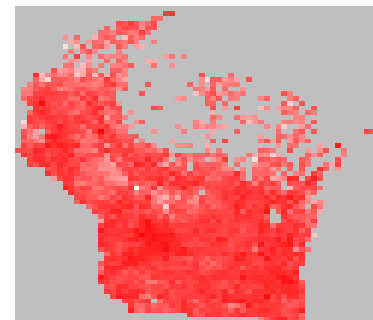
- Produce a model of per-species biomass at time of settlement
- Complicated by the presence of zeros

Most common taxa

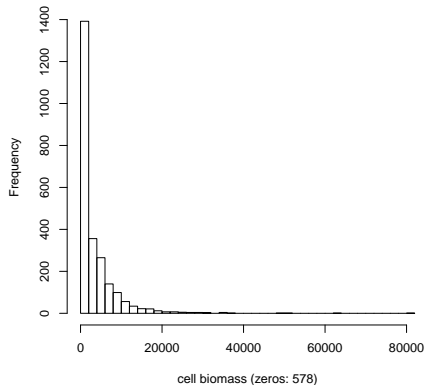
Taxon	Biomass
Oaks	7900000
Pine	6900000
Hemlock	6500000
Birches	6400000
Maple	4900000
Basswood	1700000
Elms	1600000
Tamarack	1600000
Cedar	1500000

Oaks

Observed Oaks biomass (log scale)

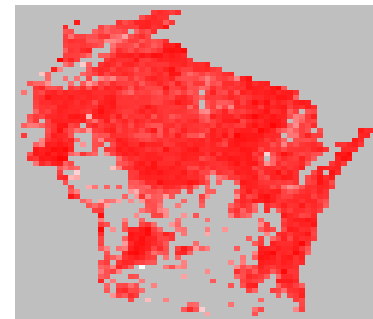


Oaks biomass

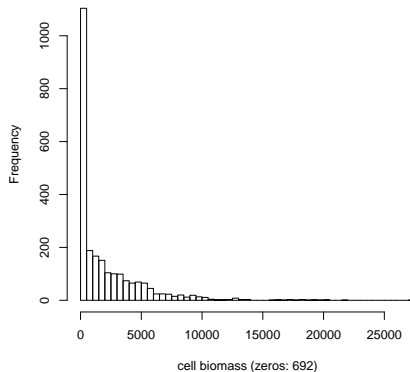


Pine

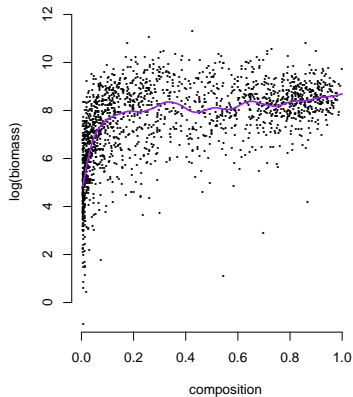
Observed Maple biomass (log scale)



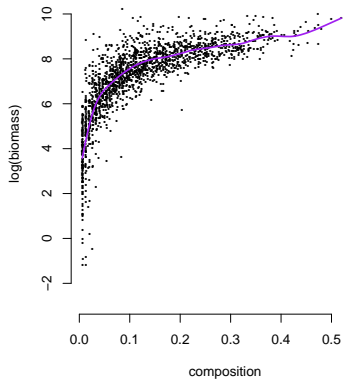
Maple biomass



Oaks biomass vs. composition



Maple biomass vs. composition



Modeling grid

Model	Spatial fit		
	<i>Indep.</i>	<i>Splines</i>	<i>GMRF</i>
One-stage (Tweedie)			
Two-stage (Bernoulli-Gamma)			

- Model types:
 - ▶ Tweedie: single stage model that accounts for exact zeros
 - ▶ Binomial-Gamma: binomial stage for presence/absence and gamma stage for biomass, conditional on presence
- Spatial fitting methods:
 - ▶ Independent: ignore spatial structure, treat observations as conditionally independent
 - ▶ Splines: smoothing spline for the spatial effect
 - ▶ GMRF: spatial random effect (fit with the INLA algorithm)

Defining terms

Let:

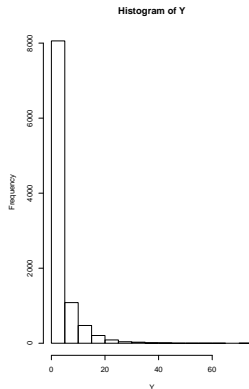
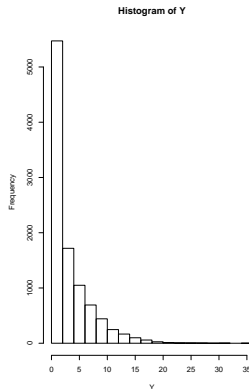
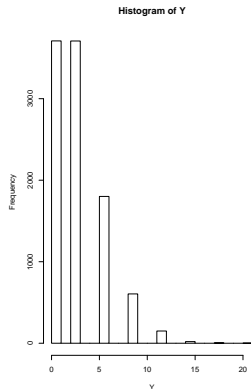
- s index location, k index taxon
- $Y_{k,s}$ denote the biomass of taxon k in cell s
- $p_{k,s}$ denote the composition fraction of taxon k in cell s
- γ_s denote the overall stem density in grid cell s

Two-stage models

- First stage: zero/non-zero
 - ▶ Logistic regression
 - ▶ $Z_s \sim \text{Bernoulli}(\zeta_s)$
 - ▶ $\text{logit}(\zeta_s) = f(\dots)$
- Second stage: distribution of positive biomass
 - ▶ $Y_s | Z_s = 1 \sim \text{Gamma}(\alpha_s, \beta_s)$
 - ▶ $E(Y_s | Z_s = 1) = \mu_s = g(\eta_s)$
 - ▶ $\eta_s = h(\dots)$

One-stage models: the Tweedie family

- Tweedie-family distributions have a point mass at zero as well as a continuous positive distribution.
- Parameter θ controls the mixture, from $\theta = 1$ (Poisson) to $\theta = 2$ (Gamma)



Conceptualizing the Tweedie distribution as a Gamma-Poisson mixture with parameters α , β , and λ :

- Draw $N \sim \text{Poisson}(\lambda)$
- Now make N iid draws: $V_\ell \sim \text{Gamma}(\alpha, \beta)$
- $Y = \sum_{\ell=1}^N V_\ell$

Tweedie model

With θ given, the Tweedie distribution is in the exponential family.

- $EY = \mu$
- $\text{var}(Y) = \phi\mu^\theta$
- ϕ is a scale parameter
- $P(Y = 0) = \exp\left(-\phi^{-1}\frac{\mu^{2-\theta}}{2-\theta}\right)$
 - ▶ $P(Y = 0) \uparrow$ as $\mu \rightarrow -\infty$
 - ▶ $P(Y = 0) \uparrow$ as $\phi \uparrow$
 - ▶ $P(Y = 0) \uparrow$ as $\theta \rightarrow 1$
- But we first need θ
 - ▶ Find $\hat{\theta}$ so that the model's deviance residuals match the assumed variance function.

Independent grid cells

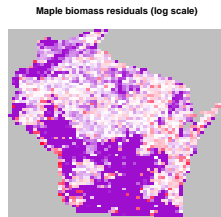
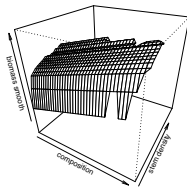
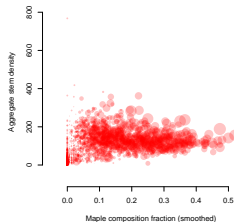
Ignoring the spatial structure, model the biomass as a function of composition and stem density

$$\eta_s = f(p_{k,s}, \gamma_s)$$

$$\mu_s = EY_s = \exp(\eta_s)$$

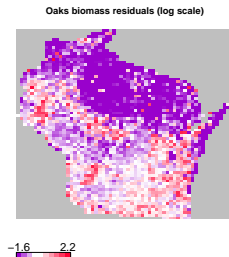
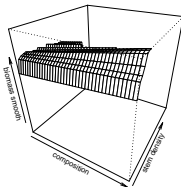
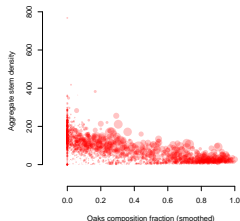
$$Y_s \sim \text{Gamma}(\alpha, \beta)$$

The pictures are from a one-stage model for Maple:



Independent grid cells

Pictures here are from a one-stage model for Oak.



Spline models

One way to account for spatial patterns is by using splines to model a spatial effect. Essentially the splines are a smooth function of latitude and longitude, and are fit using the same software as was used for the smooth functions of composition and stem density in the independent grid cell models.

Hierarchical model

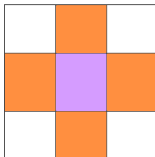
Consider a log-normal model where the mean is a Gaussian Process:

$$\begin{aligned}\log(Y_s) &\sim \mathcal{N}(g(s), \sigma^2) \\ g(\cdot) &\sim \text{GP}(f(\cdot), \Sigma)\end{aligned}$$

Parameters for such a Gaussian Process are often estimated by approximating it as a Gaussian Markov Random Field (GMRF)

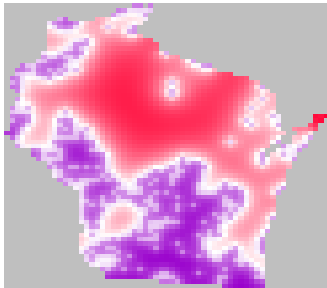
Gaussian Markov Random Field

- Markov property says that all the relevant information for modeling a grid cell is found in its neighbors.
- $g_i | \mathbf{g}_{-i}, \kappa \sim \mathcal{N} \left(n_i^{-1} \sum_{j \in \partial_i} g_j, (n_i \kappa)^{-1} \right)$

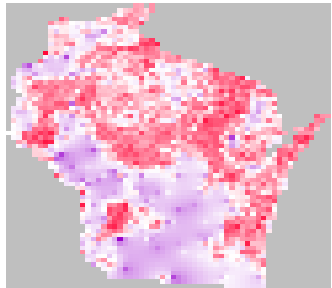


The method of integrated Nested Laplace Approximations (INLA) can be used for parameter estimation, is much faster than MCMC. The INLA software can be used for two-stage models, but the Tweedie likelihood is not currently included.

Maple logit(probability of nonzero biomass) (first stage)



Maple fitted biomass (log scale) (second stage)



Modeling grid

Model	Spatial fit		
	<i>Indep.</i>	<i>Splines</i>	<i>GRF</i>
One-stage (Tweedie)	✓	✓	
Two-stage (Bernoulli-Gamma)	✓	✓	✓

COZIGAM: an alternative for fitting the two-stage model

- COZIGAM stands for constrained zero-inflated GAM
- “Constrained” means that the latent predictors for the $P(y = 0)$ part and the $f(y|y > 0)$ part must be proportional
- Package was removed from CRAN last year
- In testing with the archived version, algorithm convergence was elusive