# Model Choice and Variable Selection in Geoadditive Regression Models

Thomas Kneib

Department of Statistics
Ludwig-Maximilians-University Munich

joint work with

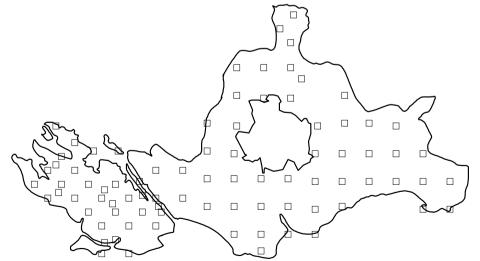Torsten Hothorn                                    Gerhard Tutz

11.2.2008

# Geoadditive Regression: Forest Health Example

- Aim of the study: Identify factors influencing the health status of trees.

- Database: Yearly visual forest health inventories carried out from 1983 to 2004 in a northern Bavarian forest district.

- 83 observation plots of beeches within a 15 km times 10 km area.

- Response: binary defoliation indicator $y_{it}$ of plot $i$ in year $t$
  ($1 =$ defoliation higher than 25%).

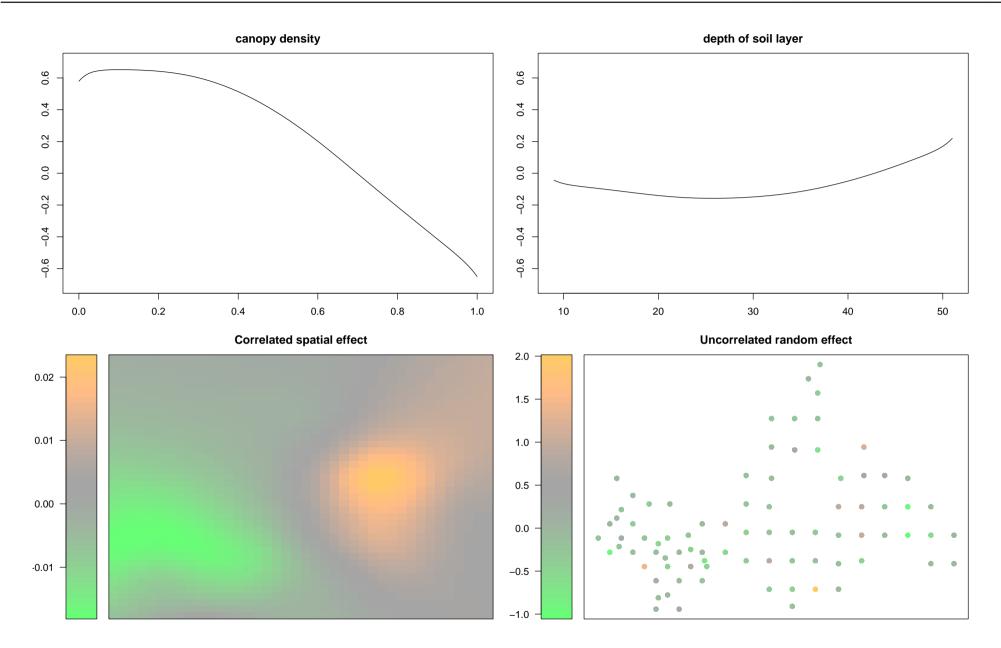- Spatially structured longitudinal data.

- **Covariates**:

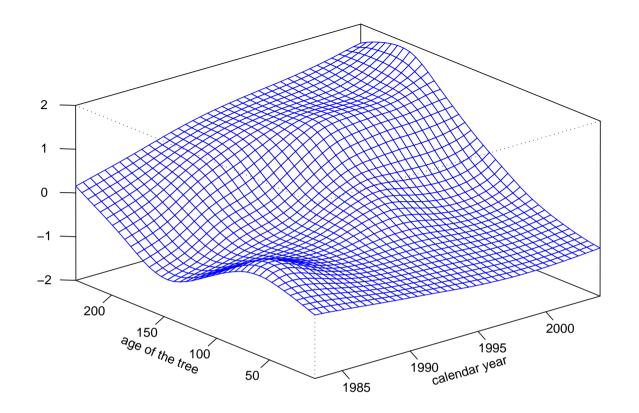  | Continuous: | average age of trees at the observation plot |
  | | elevation above sea level in meters |
  | | inclination of slope in percent |
  | | depth of soil layer in centimeters |
  | | pH-value in $0 - 2$cm depth |
  | | density of forest canopy in percent |
  | Categorical | thickness of humus layer in 5 ordered categories |
  | | level of soil moisture |
  | | base saturation in 4 ordered categories |
  | Binary | type of stand |
  | | application of fertilisation |

- Possible model:

$$P(y_{it} = 1) = \frac{\exp(\eta_{it})}{1 + \exp(\eta_{it})}$$

where $\eta_{it}$ is a <span style="color:red">geoadditive predictor</span> of the form

$$
\begin{aligned}
\eta_{it} \quad = \quad & f_1(\text{age}_{it}, t) + && \text{\color{red}interaction} \text{ between age and calendar time.} \\
& f_2(\text{canopy}_{it}) + && \text{\color{red}smooth effects} \text{ of the canopy density and} \\
& f_3(\text{soil}_{it}) + && \text{the depth of the soil layer.} \\
& f_{spat}(s_{ix}, s_{iy}) + && \text{structured and} \\
& b_i + && \text{unstructured \color{red}spatial random effects.} \\
& x'_{it}\beta && \text{parametric effects of type of stand, fertilisation,} \\
& && \text{thickness of humus layer, level of soil moisture} \\
& && \text{and base saturation.}
\end{aligned}
$$

- Questions:

  – How do we estimate the model? ⇒ Inference

  – How do we come up with the model specification? ⇒ Model choice and variable selection

⇒ Componentwise boosting for geoadditive regression models.

# Base-Learners For Geoadditive Regression Models

- Base-learning procedures for geoadditive regression models can be derived from univariate Gaussian smoothing approaches, e.g.

$$y = g(x) + \varepsilon$$

  where $\varepsilon \sim N(0, \sigma^2 I)$ and $g(x)$ is smooth.

- Spline smoothing: Approximate a function $g(x)$ by a linear combination of B-spline basis functions, i.e.
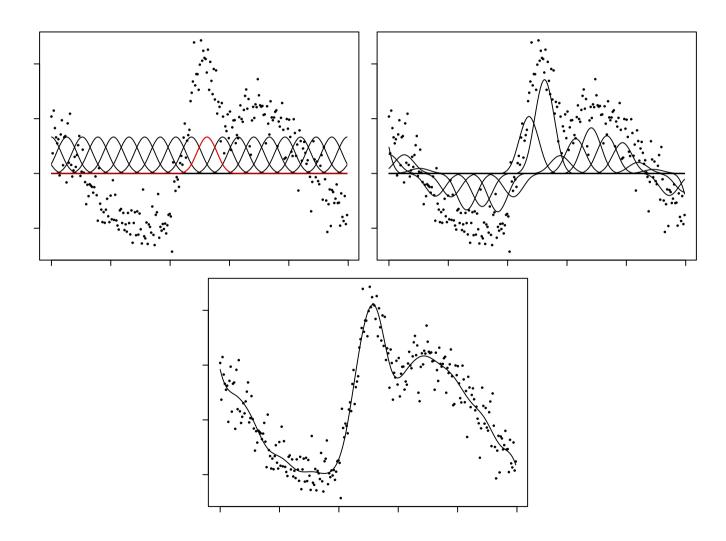
$$g(x) = \sum_j \beta_j B_j(x)$$

- In matrix notation:

$$y = X\beta + \varepsilon.$$

- Least squares estimate for $\beta$ and predicted values:

$$\hat{\beta} = (X'X)^{-1}X'y \qquad \hat{y} = X'(X'X)^{-1}X'y$$

- B-spline fit depends on the number and location of basis functions

  $\Rightarrow$ Difficult to obtain a suitable compromise between smoothness and fidelity to the data.

- Add a roughness penalty term to the least squares criterion.

- Simple approximation to squared derivative penalties: Difference penalties

$$\text{pen}(\beta) = \lambda \sum_j (\beta_j - \beta_{j-1})^2 \quad \text{or} \quad \text{pen}(\beta) = \lambda \sum_j (\beta_j - 2\beta_{j-1} + \beta_{j-2})^2.$$

- Can be written as quadratic forms

$$\lambda \beta' D' D \beta = \lambda \beta' K \beta$$

  based on difference matrices $D$.

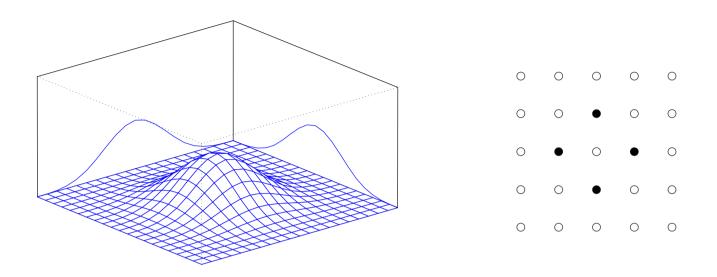- Replace the least-squares estimate and fit with penalised least squares (PLS) variants:

$$\hat{\beta} = (X'X + \lambda K)^{-1}X'y \qquad \hat{y} = X'(X'X + \lambda K)^{-1}X'y$$

- The base-learner is characterised by the hat matrix

$$S_\lambda = X'(X'X + \lambda K)^{-1}X'.$$

- PLS base-learners can also be derived for

  - Interaction surfaces $f(x_1, x_2)$ and spatial effects $f(s_x, s_y)$,

  - Varying coefficient terms $x_1 f(x_2)$ or $x_1 f(s_x, s_y)$,

  - Random intercepts $b_i$ and random slopes $x b_i$, and

  - Fixed effects $x\beta$.

- Interaction surfaces $f(x_1, x_2)$ and spatial effects $f(s_x, s_y)$:



- Define bivariate <span style="color:red">Tensor product</span> basis functions

$$B_{jk}(x_1, x_2) = B_j(x_1)B_k(x_2).$$

- Based on penalty matrices $K_1$ and $K_2$ for univariate fits define <span style="color:red">rowwise and columnwise penalties</span> as

$$
\begin{aligned}
\mathrm{pen}_{\mathsf{row}}(\beta) &= \lambda \beta' \left( I \otimes K_1 \right) \beta \\
\mathrm{pen}_{\mathsf{col}}(\beta) &= \lambda \beta' \left( K_2 \otimes I \right) \beta.
\end{aligned}
$$

- The overall penalty is then given by

$$
\mathrm{pen}(\beta) = \lambda \beta' \underbrace{\left( I \otimes K_1 + K_1 \otimes I \right)}_{} = K\beta.
$$

- Varying coefficient terms $x_1 f(x_2)$ or $x_1 f(s_x, s_y)$:

$$X = \mathrm{diag}(x_{11}, \ldots, x_{n1})X^*$$

  where $X^*$ is the design matrix representing $f(x_2)$ or $f(s_x, s_y)$.

- Cluster-specific random intercepts: The design matrix is a zero/one incidence matrix linking observations to clusters and the penalty matrix is a diagonal matrix.

- Fixed effects: Set the smoothing parameter to zero (unpenalised least squares fit).

- All base-learners can be described in terms of a <span style="color:red">penalised hat matrix</span>

$$S_\lambda = X'(X'X + \lambda K)^{-1}X'$$

  with suitably chosen design matrix $X$ and penalty matrix $K$.

# Complexity Adjustment

- The flexibility of penalised least squares base-learners depends on the choice of the smoothing parameter.

- Typical strategy: fix the smoothing parameter at a large pre-specified value.

- Difficult when comparing fixed effects, nonparametric effects and spatial effects.

    $\Rightarrow$ More flexible base-learners will be preferred in the boosting iterations leading to potential selection (and estimation) bias.

- We need an intuitive measure of complexity.

- The complexity of a linear model can be assessed by the trace of the hat matrix, since

$$\mathrm{trace}(X(X'X)^{-1}X') = \mathrm{ncol}(X).$$

- In analogy, the effective degrees of freedom of a penalised least-squares base-learner are given by

$$\mathrm{df}(\lambda_j) = \mathrm{trace}(X_j(X_j'X_j + \lambda_j K_j)^{-1}X_j').$$

- Choose the smoothing parameters for the base-learners such that

$$\mathrm{df}(\lambda_j) = 1.$$

- Difficulty: For most PLS base-learners, the penalty matrix $K$ has a non-trivial null space, i.e.

$$\dim(\mathcal{N}(K)) \geq 1.$$

- For example, a polynomial of order $k-1$ remains unpenalised for penalised splines with $k$-th order difference penalty.

$\Rightarrow \mathrm{df}(\lambda_j) = 1$ can not be achieved.

- A reparameterisation has to be applied, leading for example to

$$f(x) = \beta_0 + \beta_1 x + \ldots + \beta_{k-1} x^{k-1} + f_{\text{centered}}(x).$$

- Assign separate base-learners to the parametric components and a one degree of freedom PLS base-learner to the centered effect.

- This will also allow to choose between linear and nonlinear effects within the boosting algorithm.

# A Generic Boosting Algorithm

- Generic representation of geoadditive models:

$$\eta(\cdot) = \beta_0 + \sum_{j=1}^{r} f_j(\cdot)$$

  where the functions $f_j(\cdot)$ represent the candidate functions of the predictor.

- Componentwise boosting procedure based on the loss function $\varrho(\cdot)$:

1. Initialize the model components as $\hat{f}_j^{[0]}(\cdot) \equiv 0$, $j = 1, \ldots, r$. Set the iteration index to $m = 0$.

2. Increase $m$ by 1. Compute the current negative gradient

$$u_i = -\left.\frac{\partial}{\partial \eta}\varrho(y_i, \eta)\right|_{\eta = \hat{\eta}^{[m-1]}(\cdot)}, \quad i = 1, \ldots, n.$$

3. Choose the base-learner $g_{j*}$ that minimizes the $L_2$-loss, i.e. the best-fitting function according to

$$j^* = \underset{1 \leq j \leq r}{\operatorname{argmin}} \sum_{i=1}^{n} (u_i - \hat{g}_j(\cdot))^2$$

   where $\hat{g}_j = S_j u$.

4. Update the corresponding function estimate to

$$\hat{f}_{j*}^{[m]}(\cdot) = \hat{f}_{j*}^{[m-1]}(\cdot) + \nu S_{j*} u,$$

   where $\nu \in (0, 1]$ is a step size. For all remaining functions set

$$\hat{f}_{j}^{[m]}(\cdot) = \hat{f}_{j}^{[m-1]}(\cdot), \quad j \neq j^*.$$

5. Iterate steps 2 to 4 until $m = m_{\text{stop}}$.

- Determine $m_{\text{stop}}$ based on AIC reduction or cross-validation.

- Boosting implements both variable selection and model choice:

  - Variable selection: Stop the boosting procedure after an appropriate number of iterations (for example based on AIC reduction).

  - Model choice: Consider concurring base-learning procedures for the same covariate, e.g. linear vs. nonlinear modeling.
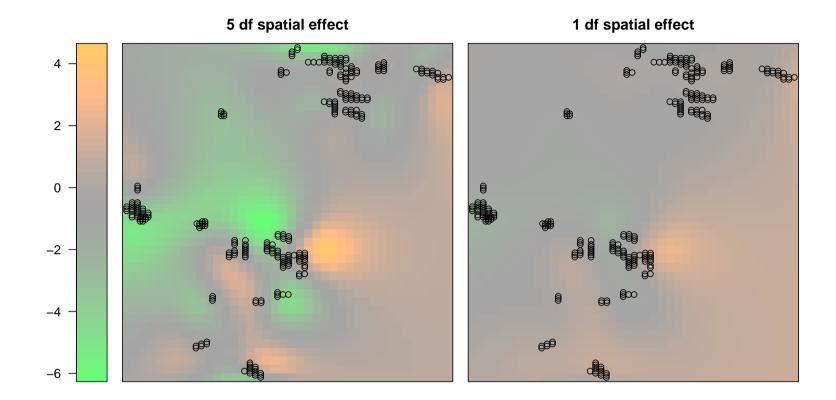
# Habitat Suitability Analyses

- Identify factors influencing habitat suitability for breeding bird communities collected in seven structural guilds (SG).

- Variable of interest: Counts of subjects from a specific structural guild collected at 258 observation plots in a Northern Bavarian forest district.

- Research questions:

a) Which covariates influence habitat suitability (31 covariates in total)? Does spatial correlation have an impact on variable selection?

b) Are there nonlinear effects of some of the covariates?

c) Are effects varying spatially?

- All questions can be addressed with the boosting approach.

# Variable Selection in the presence of spatial correlation

- Selection frequencies in a spatial Poisson-GLM:

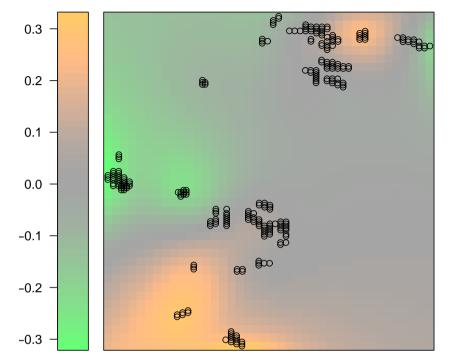| | GST | DBH | AOT | AFS | DWC | LOG | SNA | COO |
|---|---|---|---|---|---|---|---|---|
| non-spatial GLM | 0 | 0 | 0 | 0.06 | 0.3 | 0 | 0.01 | 0 |
| spatial with 5 df | 0 | 0.02 | 0 | 0.01 | 0.05 | 0 | 0.01 | 0 |
| spatial with 1 df | 0 | 0 | 0 | 0.06 | 0.15 | 0 | 0 | 0 |
| | COM | CRS | HRS | OAK | COT | PIO | ALA | MAT |
| non-spatial GLM | 0.03 | 0.04 | 0.03 | 0.05 | 0.06 | 0 | 0.04 | 0.06 |
| spatial with 5 df | 0 | 0.01 | 0 | 0 | 0 | 0 | 0.01 | 0.05 |
| spatial with 1 df | 0.03 | 0.02 | 0.02 | 0.04 | 0.05 | 0 | 0.03 | 0.04 |
| | GAP | AGR | ROA | LCA | SCA | HOT | CTR | RLL |
| non-spatial GLM | 0.03 | 0 | 0 | 0.1 | 0.07 | 0 | 0 | 0 |
| spatial with 5 df | 0.01 | 0 | 0.01 | 0.01 | 0.01 | 0 | 0 | 0 |
| spatial with 1 df | 0.03 | 0 | 0 | 0.07 | 0.06 | 0 | 0 | 0 |
| | BOL | MSP | MDT | MAD | COL | AGL | SUL | spatial |
| non-spatial GLM | 0 | 0.06 | 0 | 0 | 0.05 | 0 | 0 | 0 |
| spatial with 5 df | 0 | 0 | 0 | 0 | 0.03 | 0 | 0 | 0.76 |
| spatial with 1 df | 0 | 0.04 | 0 | 0 | 0.04 | 0 | 0 | 0.3 |

- Spatial effects for high and low degrees of freedom (SG4):
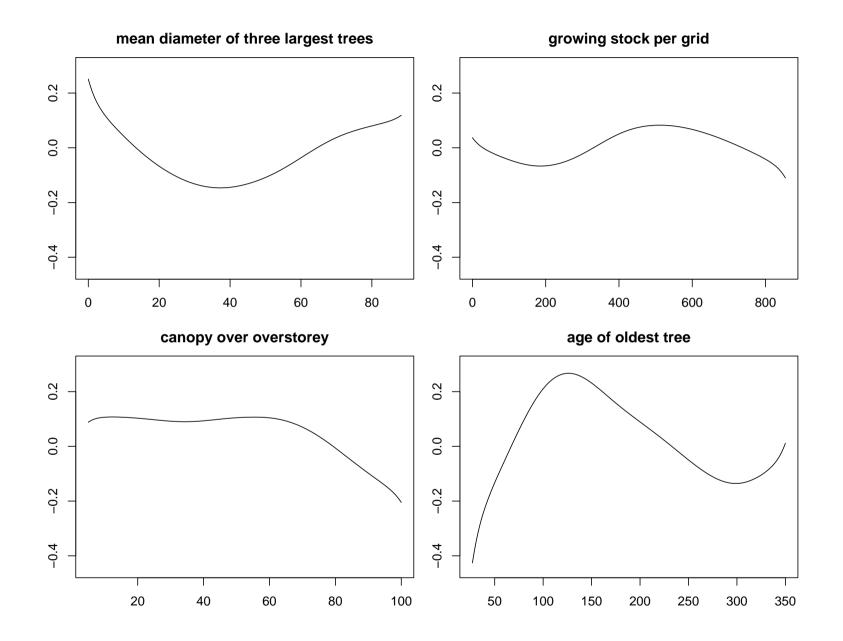


- Spatial correlation has non-negligible influence on variable selection.

- Making terms comparable in terms of complexity is essential to obtain valid results.

# Geoadditive Models

- Instead of linear modelling, allow for nonlinear effects of all 31 covariates.

- Decompose nonlinear effects into a linear part and a nonlinear part with one degree of freedom.

- Variable selection for SG5 results in 7 variables without any influence, 3 linear effects, and 21 nonlinear effects.
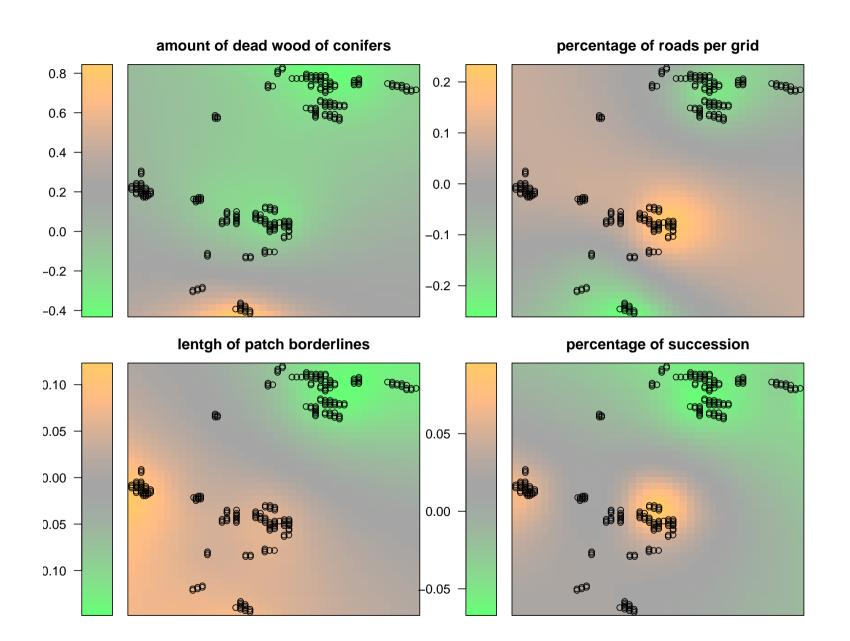
**SG5: Geoadditive Model**

# Space-varying effects

- Instead of allowing for nonlinear effects, consider space-varying effects $xg(s_x, s_y)$ for all covariates.

- Decompose space-varying effects into a linear part and a space-varying part with one degree of freedom.

- For SG3, 6 variables have no influence at all, 13 variables have linear effects, and 12 variables are associated with space-varying effects.

- The spatial effect is completely explained by the space-varying effects of the covariates.

# Summary & Extensions

- Generic boosting algorithm for model choice and variable selection in geoadditive regression models.

- Avoid selection bias by careful parameterisation.

- Implemented in the R-package **mboost**.

- Future plans:

  - Derive base-learning procedures for other types of spatial effects (regional data, anisotropic spatial effects).

  - Construct spatio-temporal base-learners based on tensor product approaches.

  - Extend methodology to model selection in continuous time survival models.

- Reference: Kneib, T., Hothorn, T. and Tutz, G.: Model Choice and Variable Selection in Geoadditive Regression. Under revision for *Biometrics*.

- Find out more:

$$\texttt{http://www.stat.uni-muenchen.de/\~{}kneib}$$