

Structured additive regression for multitegogical space-time data: a mixed model approach

Thomas Kneib

Department of Statistics, University of Munich

1. Forest health data
2. Regression models for ordinal responses
3. Structured additive regression
4. Mixed model representation
5. Results
6. Software
7. Discussion

Forest health data

- Data collected in yearly forest health inventories carried out in a forest in northern Bavaria from 1983 to 2001.
- 83 observation points with beeches in an area extending 15 km from east to west and 10 km from north to south.
- y_{it} , the defoliation degree of beech i in year t , is measured in three **ordered categories** (multicategorical response):
 - $y_{it} = 1$ no defoliation,
 - $y_{it} = 2$ defoliation 25% or less,
 - $y_{it} = 3$ defoliation above 25%.
- Covariates:
 - t calendar time,
 - s_i site of the beech,
 - a_{it} age of the tree in years,
 - u_{it} further (mostly categorical) covariates.

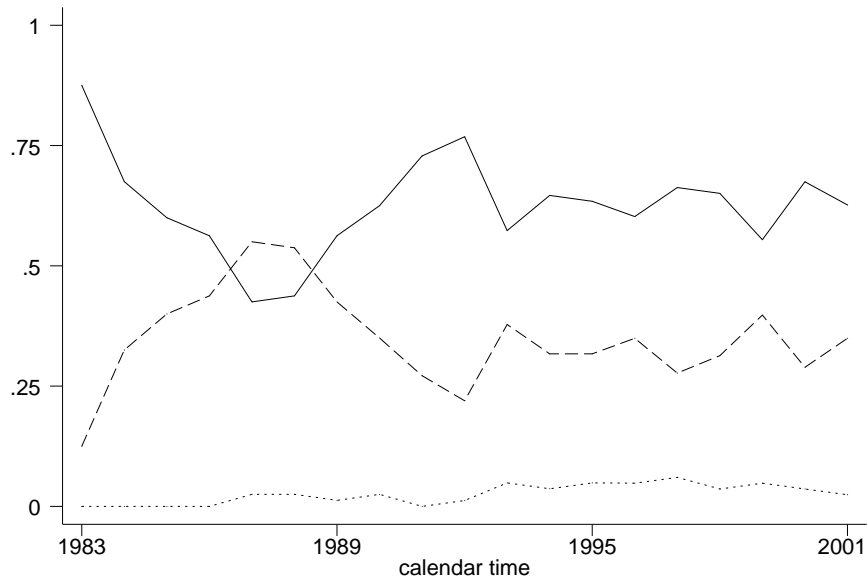
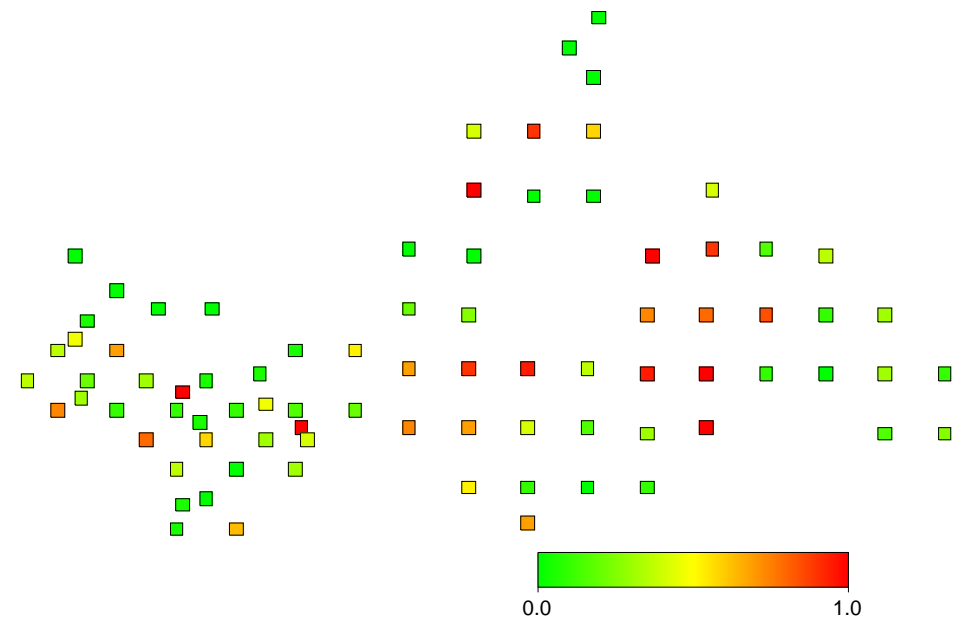


Figure 1: Temporal development of the frequency of the damage states:

- no damage,
- - - medium damage,
- ... severe damage.

Figure 2: Spatial distribution of the beeches and percentage of time points for which a beech was classified to be damaged (damage state 2 or 3).



Regression models for ordinal responses

- Response y_{it} follows multinomial distribution with three ordered categories $r = 1, 2, 3$.
- Model the **cumulative probabilities**

$$P(y_{it} \leq r) = F(\theta_r - \eta_{it})$$

with thresholds $-\infty = \theta_0 < \theta_1 < \theta_2 < \theta_3 = \infty$ and linear predictor η_{it} .

- $F(\cdot)$ can be any cumulative distribution function:
 - standard normal \implies cumulative **probit** model,
 - logistic \implies cumulative **logit** model.

- Consider a random variable with density $f = F'$ and expectation η_{it} .

⇒ Linear predictor determines **shift on latent scale**.

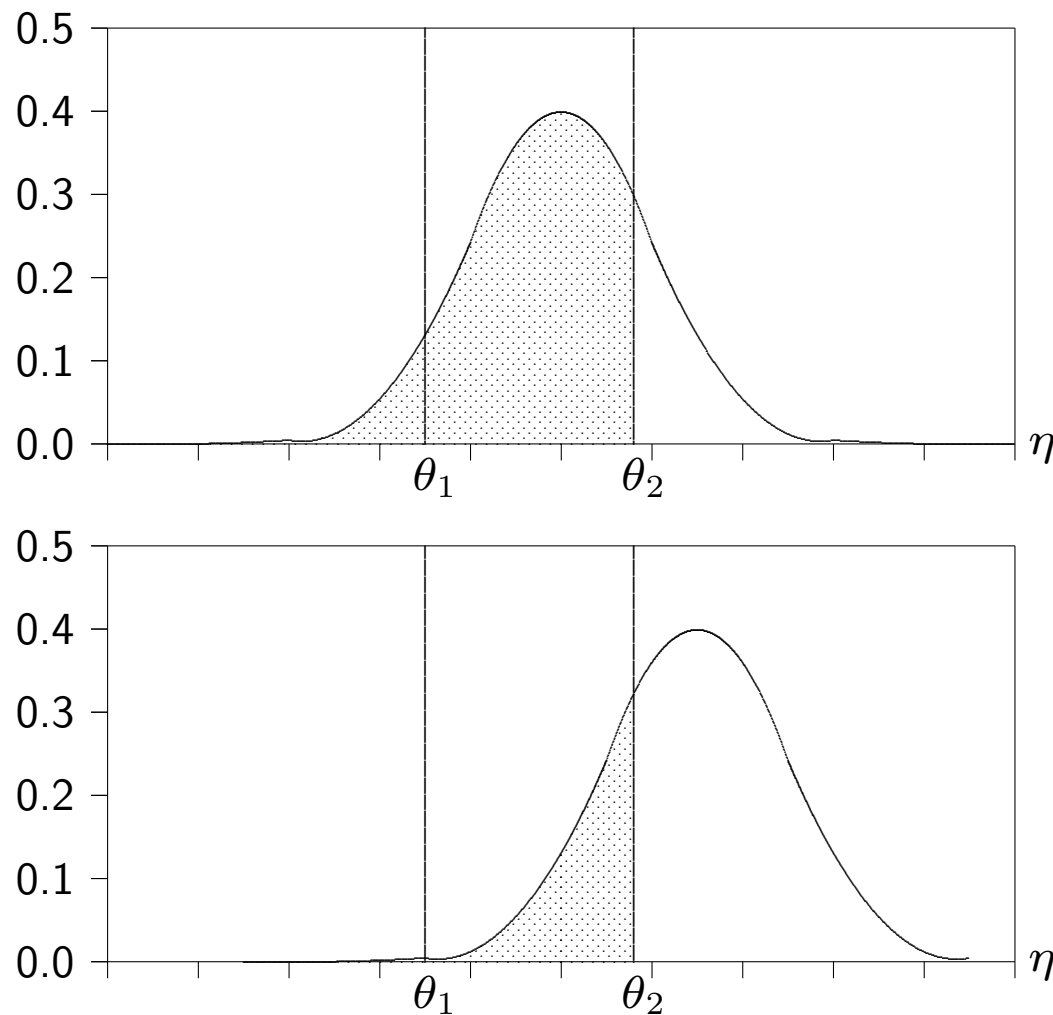


Figure 3: The shaded areas represent $P(y_{it} \leq 2)$ for different values of η_{it} .

- Limitations of a purely parametric approach:
 - Spatio-temporal structure of the data implies **spatial** and **temporal correlations**.
 - **Nonlinear effects** of continuous covariates?
 - **Complex interactions** between covariates?
- ⇒ Structured additive regression models.

Structured additive regression

- Replace usual parametric predictor with a **flexible semiparametric** predictor

$$\eta_{it} = f_1(t) + f_2(a_{it}) + f_3(t, a_{it}) + f_{spat}(s_i) + u'_{it}\gamma,$$

where

- f_1 and f_2 are **nonparametric** functions of calendar time and age,
 - f_3 is an **interaction surface** between calendar time and age,
 - f_{spat} is a **spatial** function, and
 - u is a vector of further covariates with parametric effects.
- Structured additive regression extends (and combines) generalized additive mixed models, geoadditive models and varying coefficient models.
 - Allows **unified treatment** of all effects within a **Bayesian framework**.

- $f_1(t), f_2(a_{it})$: **P-splines**
 - Approximate f_j by a B-spline of a certain degree (basis function approach).
 - Penalize differences between parameters of adjacent basis functions to ensure smoothness.
 - Alternatives: **Random walks**, more general **autoregressive priors**.
- $f_3(t, a_{it})$: **Two-dimensional extensions of P-splines**
 - Define two-dimensional basis functions based on tensor products of one-dimensional B-splines.
 - Use priors from spatial statistics for penalization.
 - Alternative: **Varying coefficient models**, if one of the interacting variables is categorical.

- $f_{spat}(s_i)$: **Markov random fields**
 - Consider two trees as neighbors if their distance is less than (e.g.) 1.2 km.
 - Assume that the expected value of $f_{spat}(s)$ is the average of the function evaluations of adjacent sites.
- $f_{spat}(s_i)$: **Stationary Gaussian random fields** (kriging)
 - Spatial effect follows a zero mean stationary Gaussian stochastic process.
 - Correlation of two arbitrary sites is defined by an intrinsic correlation function.
- Split up spatial effect into **structured** and **unstructured** part:

$$f_{spat}(s_i) = f_{str}(s_i) + f_{unstr}(s_i)$$

The unstructured effect can be modelled by i.i.d. random effects, the structured effect by a MRF or a GRF.

- All effects f_j can be expressed as the product of a **design matrix** Z_j and a vector of **regression coefficients** β_j .
- Rewrite the structured additive predictor in matrix notation as

$$\eta = Z_1\beta_1 + Z_2\beta_2 + Z_3\beta_3 + Z_{spat}\beta_{spat} + U\gamma.$$

- Bayesian approach: Assign an appropriate **prior** to β_j .
- All priors can be cast into the **general form**

$$p(\beta_j|\tau_j^2) \propto \exp\left(-\frac{1}{2\tau_j^2}\beta_j'K_j\beta_j\right)$$

where K_j is a **penalty matrix** and τ_j^2 is a **smoothing parameter**.

- Type of the covariate and prior beliefs about the smoothness of f_j determine special Z_j and K_j .

Mixed model representation

- Each parameter vector β_j can be partitioned into an **unpenalized part** (with flat prior) and a **penalized part** (with i.i.d. Gaussian prior) yielding a **variance components model**

$$\eta = X^{unp} \beta^{unp} + X^{pen} \beta^{pen}$$

with

$$p(\beta^{unp}) \propto \text{const} \quad \beta^{pen} \sim N(0, \Lambda)$$

and

$$\Lambda = \text{blockdiag}(\tau_1^2 I, \dots, \tau_4^2 I).$$

- Regression coefficients are estimated via **modified Fisher scoring**.
- The mixed model representation allows for **restricted maximum likelihood** / **marginal likelihood** estimation of the variance components:

$$L(\Lambda) = \int L(\beta^{unp}, \beta^{pen}, \Lambda) p(\beta^{pen}) d\beta^{pen} d\beta^{unp} \rightarrow \max_{\Lambda}.$$

- From a Bayesian perspective, we get **empirical Bayes** / **posterior mode** estimates.
- Closely related to penalized likelihood.
- Fahrmeir, Kneib and Lang (2004) derive numerically efficient formulae that allow the computation even for fairly large data sets.

Results

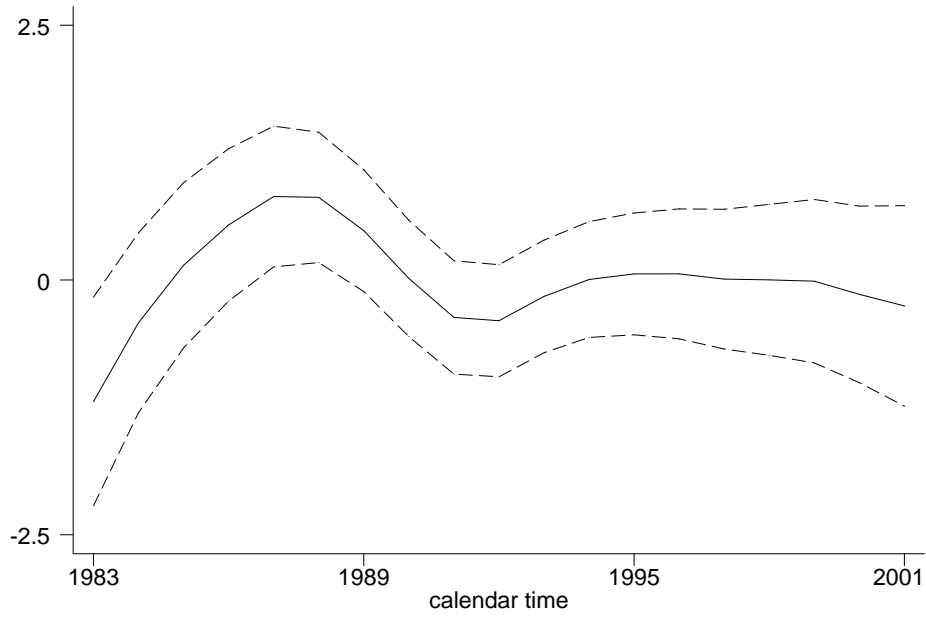


Figure 4: Time trend.

Figure 5: Age effect.

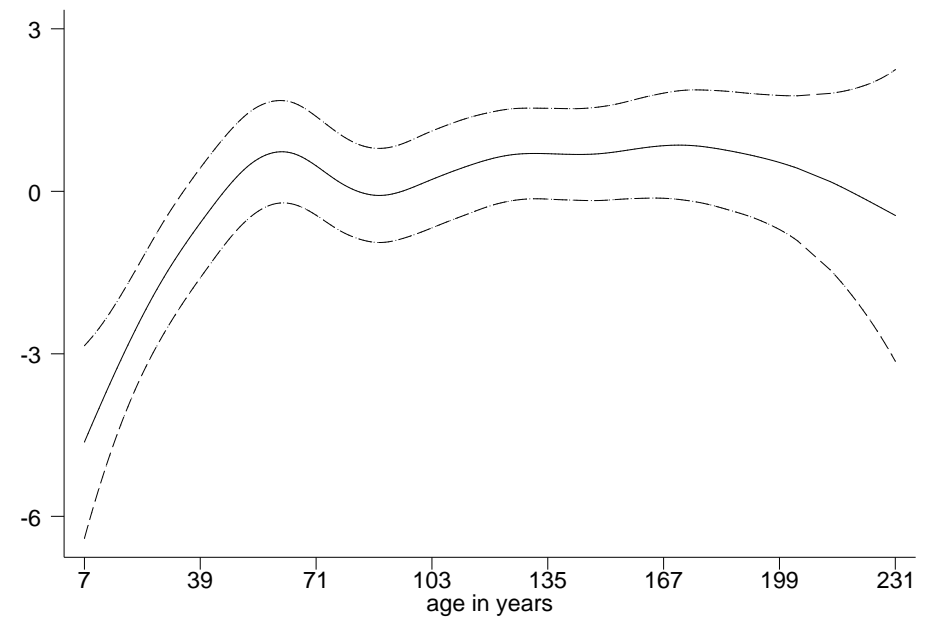


Figure 6: Structured spatial effect.

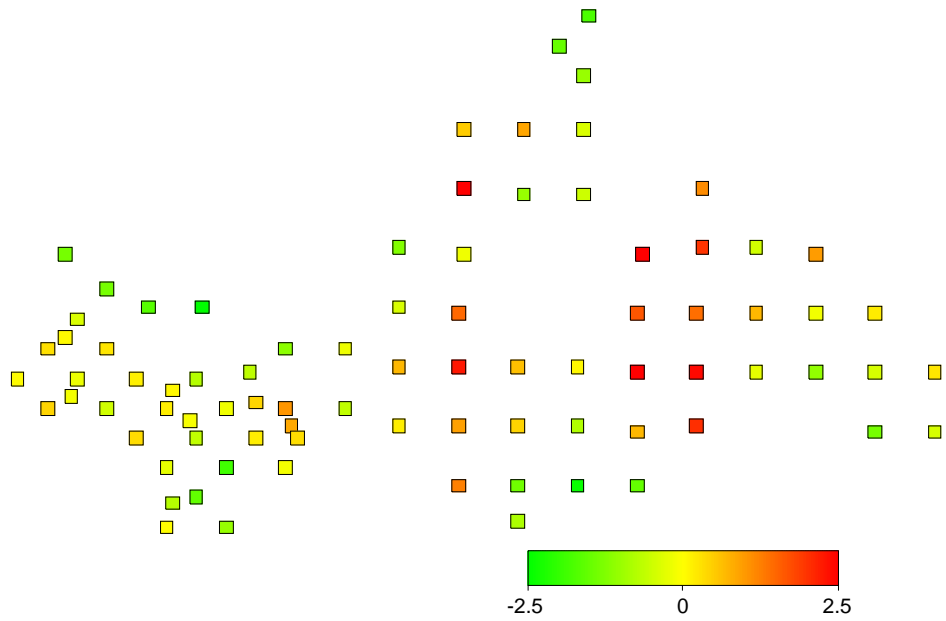
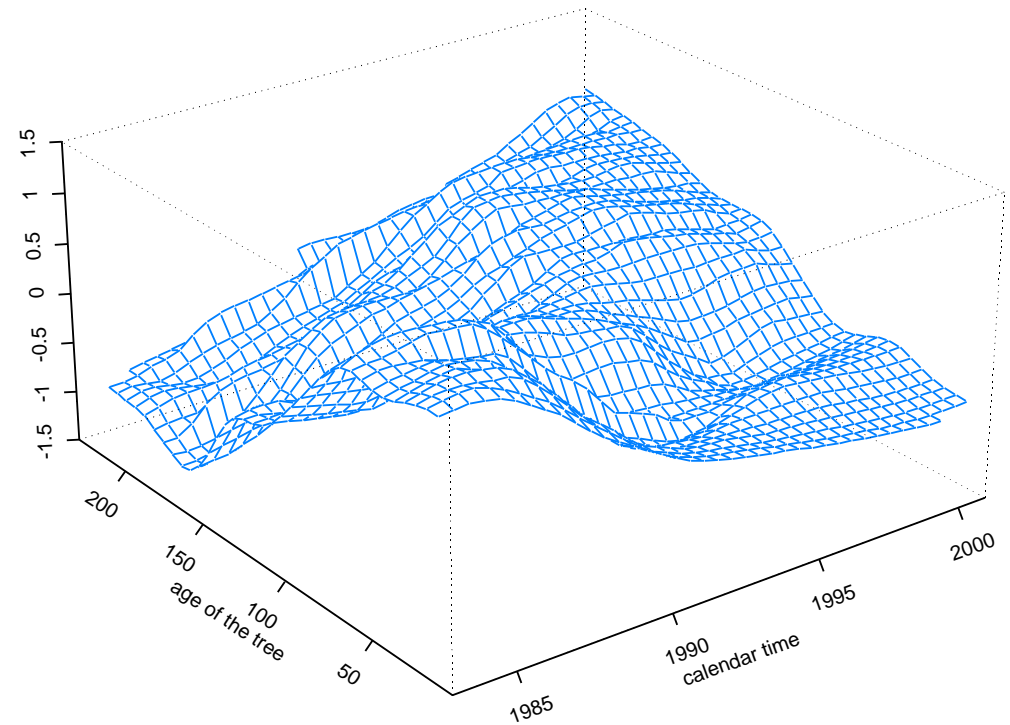


Figure 7: Interaction effect.



Software

- Estimation was carried out using BayesX, a public domain software package for Bayesian inference.



- Available from

<http://www.stat.uni-muenchen.de/~lang/bayesx>

- Features (within a mixed model setting):
 - Responses: Gaussian, Gamma, Poisson, Binomial, ordered and unordered multinomial.
 - Continuous covariates and time scales: Random Walks, P-splines, autoregressive priors for seasonal components.
 - Spatial Covariates: Markov random fields, stationary Gaussian random fields, two-dimensional P-Splines.
 - Interactions: Two-dimensional P-splines, varying coefficient models with continuous and spatial effect modifiers.
 - Random intercepts and random slopes.

Discussion

- Models for nominal responses are supported, too.
- Comparison with fully Bayesian approach based on MCMC:

Cons:

- Credible intervals rely on asymptotic normality.
- Only plug-in estimates for functionals.

Pros:

- No questions concerning mixing and convergence.
- No sensitivity with respect to prior assumptions on variance parameters.
- Somewhat better point estimates (in simulations).

- **Future work:**
 - Multinomial probit models with correlated latent utilities.
 - Category-specific covariates in nominal models.
 - Covariate-dependent thresholds in ordinal models.

References

- Fahrmeir, L., Kneib, T. and Lang, S. (2004): Penalized structured additive regression for space-time data: A Bayesian perspective. *Statistica Sinica* (to appear).
- Kneib, T. and Fahrmeir, L. (2004): Structured additive regression for multicategorical space-time data: A mixed model approach. SFB 386 Discussion Paper 377, University of Munich.
- Both available from

<http://www.stat.uni-muenchen.de/~kneib>