

Hierarchical Structured Additive Regression

Lang, Stefan¹ Umlauf, Nikolaus² Brunauer, Wolfgang³

¹²University of Innsbruck, Austria

³Immobilien Rating GmbH, Austria

July 2010

1. Example of hierarchical data structures
2. Structured additive regression models
3. Hierarchical formulation and MCMC inference
4. Alternative sampling scheme based on transformed parametrization
5. Results: Hedonic regression data for house prices

Example of hierarchical data structures

Hedonic regression data for house prices in Austria

Variable of primary interest

house price or log house price

Covariates

- Structural (physical) characteristics, like floor space area, constructional condition, age etc., and
- neighborhood (locational) characteristics, often on various levels of aggregation, like the proximity to places of work, the social composition of the neighborhood etc.

Four-level hierarchical model

$$\begin{aligned}
 \text{level 1: } \mathbf{lnp} &= f_1(\mathbf{area}) + \cdots + f_q(\mathbf{age}) + \mathbf{v}\boldsymbol{\gamma} + f_{\text{municipal}}(\mathbf{s}_1) + \boldsymbol{\varepsilon}_1 \\
 \text{level 2: } f_{\text{municipal}}(\mathbf{s}_1) &= f_{1_1}(\mathbf{purchase\ power}) + \cdots + f_{p_1}(\mathbf{level\ of\ education}) \\
 &\quad + f_{\text{district}}(\mathbf{s}_2) + \boldsymbol{\varepsilon}_2 \\
 \text{level 3: } f_{\text{district}}(\mathbf{s}_2) &= f_{1_2}(\mathbf{unemployment\ rate}) + f_{\text{county}}(\mathbf{s}_3) + \boldsymbol{\varepsilon}_3 \\
 \text{level 4: } f_{\text{county}}(\mathbf{s}_3) &= \boldsymbol{\varepsilon}_4
 \end{aligned}$$

The f 's are possibly nonlinear functions of the covariates.

This is an example of *hierarchical structured additive regression models*.

Example of hierarchical data structures

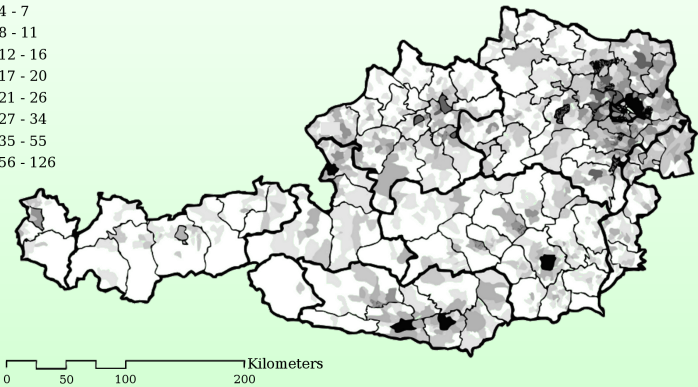
County

District

Number of observations

Missing

- 1
- 2 - 3
- 4 - 7
- 8 - 11
- 12 - 16
- 17 - 20
- 21 - 26
- 27 - 34
- 35 - 55
- 56 - 126



Structured additive regression models

- Distributional and structural assumptions, given covariates and parameters, are based on Generalized Linear Models
- $E(y|\mathbf{x}, \mathbf{v}) = h(\eta)$ with structured additive predictor

$$\eta = f_1(x_1) + \dots + f_p(x_p) + \mathbf{v}'\boldsymbol{\gamma}$$

In the following we only consider additive models with

$$y = \eta + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$

- $\mathbf{v}'\boldsymbol{\gamma}$ parametric part of the predictor
- x_j continuous covariate, time scale, location or unit-or cluster index
- x_j may be two (even higher) dimensional for modeling interactions
- f_j one-/two (even higher) dimensional, not necessarily continuous functions

Overview: Modeling the functions f_j

$f_j(x_j) = f(x)$	$x_j = x$	nonlinear effect of x
$f_j(x_j) = f_{spat}(s)$	$x_j = s$	spatial effect of location variable $\mathbf{s} = (1, 2, \dots, S)'$
$f_j(x_j) = x_2 f(x_1)$	$x_j = (x_1, x_2)$	interaction effect between x_1 and x_2
$f_j(x_j) = f_{1 2}(x_1, x_2)$	$x_j = (x_1, x_2)$	nonlinear interaction between x_1 and x_2
$f_j(x_j) = \beta_i u$	$x_j = (u, i)$	individual specific random effect with $\mathbf{u} = (1, 2, \dots, U)'$

General form

- Vector of function evaluations $\mathbf{f}_j = (f_{1j}, \dots, f_{nj})'$ can be written as:

$$\mathbf{f}_j = \mathbf{Z}_j \boldsymbol{\beta}_j$$

with \mathbf{Z}_j as the design matrix, where $\boldsymbol{\beta}_j$ are unknown regression coefficients

- Form of \mathbf{Z}_j only depends on the functional type chosen
- Penalized least squares:

$$\text{PLS}(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \|\mathbf{y} - \boldsymbol{\eta}\|^2 + \lambda_1 \boldsymbol{\beta}'_1 \mathbf{K}_1 \boldsymbol{\beta}_1 + \dots + \lambda_p \boldsymbol{\beta}'_p \mathbf{K}_p \boldsymbol{\beta}_p$$

General form

- Prior for β in the corresponding Bayesian approach

$$p(\beta_j | \tau_j^2) \propto \left(\frac{1}{2\pi\tau_j^2} \right)^{rk(\mathbf{K}_j)/2} \exp \left(-\frac{1}{2\tau_j^2} \beta_j' \mathbf{K}_j \beta_j \right) I(\mathbf{A}\beta_j = \mathbf{0})$$

τ_j^2 variance parameter, governs the smoothness of f_j , relation to frequentists by $\lambda_j = \sigma^2 / \tau_j^2$

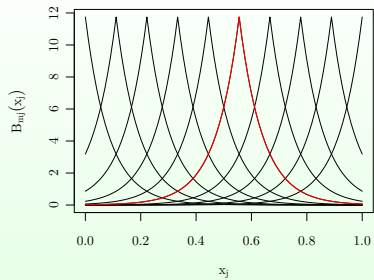
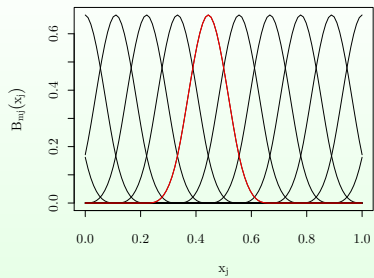
- $\mathbf{A}\beta_j = \mathbf{0}$ is an identifiability constraint, e.g. $\mathbf{A} = (1, \dots, 1)'$ such that the β 's sum up to zero
- Structure of \mathbf{K}_j also depends on the type of covariates and on assumptions about smoothness of \mathbf{f}_j

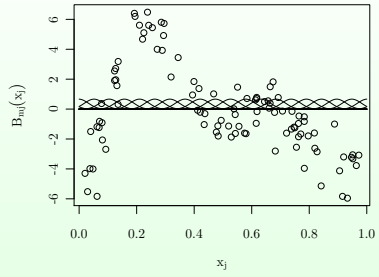
General form

- Basis functions $B_{mj}(\cdot)$ in

$$f_j(x_j) = \sum_{m=1}^{M_j} \beta_{mj} B_{mj}(x_j)$$

may include e.g. a polynomial, B-spline, Matérn basis (one or more dimensional), etc.





Hierarchical formulation and MCMC inference

Multilevel/Hierarchical structured additive model with k hierarchies within a first stage term $\mathbf{Z}_j\boldsymbol{\beta}_j$ may be written as

$$\begin{aligned}
 \mathbf{y} &= \mathbf{Z}_1\boldsymbol{\beta}_1 + \dots + \mathbf{Z}_p\boldsymbol{\beta}_p + \mathbf{v}\boldsymbol{\gamma} + \boldsymbol{\varepsilon} \\
 \boldsymbol{\beta}_j &= \mathbf{Z}_{j1_1}\boldsymbol{\beta}_{j1_1} + \dots + \mathbf{Z}_{jp_1}\boldsymbol{\beta}_{jp_1} + \mathbf{v}_j\boldsymbol{\gamma}_j + \mathbf{u}_j \\
 &\vdots \\
 \boldsymbol{\beta}_{j,j_1,\dots,j_k} &= \mathbf{z}_{j,j_1,\dots,j_k}\boldsymbol{\beta}_{j,j_1,\dots,j_k} + \dots + \mathbf{z}_{j,j_1,\dots,j_k}\boldsymbol{\beta}_{j,j_1,\dots,j_k} + \mathbf{v}_{j,j_1,\dots,j_k}\boldsymbol{\gamma}_{j,j_1,\dots,j_k} + \mathbf{u}_{j,j_1,\dots,j_k} \\
 \boldsymbol{\beta}_{j,j_1,\dots,j_k} &= \boldsymbol{\eta}_{j,j_1,\dots,j_k} + \mathbf{u}_{j,j_1,\dots,j_k}
 \end{aligned}$$

with $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{W}^{-1})$ and $\mathbf{u}_{j,j_1,\dots,j_k} \sim N(\mathbf{0}, \tau_{j,j_1,\dots,j_k}^2\mathbf{K}_{j,j_1,\dots,j_k}^{-1})$

The full conditionals for the regression coefficients are multivariate Gaussian. Starting from a first level view, the precision matrix Σ_{β_j} and mean μ_{β_j} are given by

$$\begin{aligned}\Sigma_{\beta_j} &= \sigma^2 \left(\mathbf{Z}'_j \mathbf{W} \mathbf{Z}_j + \frac{\sigma^2}{\tau_j^2} \mathbf{K}_j \right)^{-1} \\ \mu_{\beta_j} &= \Sigma_{\beta_j} \left(\frac{1}{\sigma^2} \mathbf{Z}'_j \mathbf{W} \mathbf{r} + \frac{1}{\tau_j^2} \boldsymbol{\eta}_{\beta_j} \right)\end{aligned}$$

and for the higher levels

$$\begin{aligned}\Sigma_{\beta_{j,j_1,\dots,j_k}} &= \tau_{j,j_1,\dots,j_{k-1}}^2 \left(\mathbf{Z}'_{j,j_1,\dots,j_k} \mathbf{Z}_{j,j_1,\dots,j_k} + \frac{\tau_{j,j_1,\dots,j_{k-1}}^2}{\tau_{j,j_1,\dots,j_k}^2} \mathbf{K}_{j,j_1,\dots,j_k} \right)^{-1} \\ \mu_{\beta_{j,j_1,\dots,j_k}} &= \Sigma_{\beta_{j,j_1,\dots,j_k}} \left(\frac{1}{\tau_{j,j_1,\dots,j_{k-1}}^2} \mathbf{Z}'_{j,j_1,\dots,j_k} \mathbf{r} + \frac{1}{\tau_{j,j_1,\dots,j_k}^2} \boldsymbol{\eta}_{\beta_{j,j_1,\dots,j_k}} \right)\end{aligned}$$

Properties

- *Reduced complexity in higher stages of the hierarchy:*
 - Number of “observations” in the higher levels is much less than the actual number of observations n .
 - Full conditionals for regression coefficients are Gaussian regardless of the response distribution in the first level of the hierarchy.
- *Sparsity*
Design matrices and posterior precision matrices are typically sparse (after reordering of parameters).
- *Number of different observations smaller than sample size*
Typically the number of different observations $x_{j(1_j)}, \dots, x_{j(n_j)}$ in \mathbf{Z}_j is much smaller than the total number n of observations, i.e. $n_j \ll n$.

- Denote by $z_{(1)}^{(2)} < z_{(2)}^{(2)} < \dots < z_{(m)}^{(2)}$ the m ordered different observations of $z^{(2)}$.
- Compute the index vector \mathbf{ind} with elements $\mathbf{ind}[i] \in \{1, \dots, m\}$ denoting the category of the i -th observation, i.e. if $z_i^{(2)} = z_{(j)}^{(2)}$ then $\mathbf{ind}[i] = j$.
- Decompose the design matrix in $\mathbf{Z} = \mathbf{DP}\tilde{\mathbf{Z}}$ where
 - $\tilde{\mathbf{Z}}$ is the $m \times K$ reduced design matrix for the different and sorted observations $z_{(1)}^{(2)}, \dots, z_{(m)}^{(2)}$, i.e. $\tilde{\mathbf{Z}}[s, k] = B_k \left(z_{(s)}^{(2)} \right)$, $s = 1, \dots, m$, $k = 1, \dots, K$,
 - \mathbf{P} is a $n \times m$ permutation matrix, which reverts the sorting, i.e. $\mathbf{P}[i, s] = I(\mathbf{ind}(i) = s)$.
- For the vector of function evaluations we obtain $\mathbf{f} = \mathbf{Z}\boldsymbol{\beta} = \mathbf{DP}\tilde{\mathbf{Z}}\boldsymbol{\beta}$.

We get

$$\mathbf{Z}'\mathbf{W}\mathbf{Z} = \tilde{\mathbf{Z}}'\mathbf{P}'\mathbf{D}'\mathbf{W}\mathbf{D}\mathbf{P}\tilde{\mathbf{Z}} = \tilde{\mathbf{Z}}'\tilde{\mathbf{W}}\tilde{\mathbf{Z}},$$

where

$$\tilde{\mathbf{W}} = \mathbf{P}'\mathbf{D}'\mathbf{W}\mathbf{D}\mathbf{P} = \text{diag}(\tilde{w}_1, \dots, \tilde{w}_m)$$

and the “reduced” weights \tilde{w}_s , are given by

$$\tilde{w}_s = \sum_{i: \text{ind}[i]=s} \left((z_i^{(1)})^2 \right) w_i.$$

The weights \tilde{w}_s can be computed by first initializing $\tilde{w}_s = 0$ followed by a simple loop: For $i = 1, \dots, n$ add $\left((z_i^{(1)})^2 \right) w_i$ to $\tilde{w}_{\text{ind}[i]}$.

For $\mathbf{Z}'\mathbf{W}\mathbf{r}$ we obtain

$$\mathbf{Z}'\mathbf{W}\mathbf{r} = \tilde{\mathbf{Z}}'\mathbf{P}'\mathbf{D}'\mathbf{W}\mathbf{r} = \tilde{\mathbf{Z}}'\tilde{\mathbf{r}},$$

where the $m \times 1$ vector $\tilde{\mathbf{r}} = (\tilde{r}_1, \dots, \tilde{r}_m)'$ of “reduced” partial residuals is given by

$$\tilde{r}_s = \sum_{i: \text{ind}[i]=s} z_i^{(1)} w_i r_i.$$

The \tilde{r}_s are computed by first initializing $\tilde{r}_s = 0$ followed by the loop: For $i = 1, \dots, n$ add $z_i^{(1)} w_i r_i$ to $\tilde{r}_{\text{ind}(i)}$.

Alternative sampling scheme based on transformed parametrization

- (i.) Cholesky decomposition $\mathbf{R}\mathbf{R}'$ of $\mathbf{Z}'\mathbf{W}\mathbf{Z}$
- (ii.) Singular value decomposition $\mathbf{Q}\mathbf{S}\mathbf{Q}' = \mathbf{R}^{-1}\mathbf{K}(\mathbf{R}')^{-1}$,
 $\mathbf{S} = \text{diag}(s_1, \dots, s_M)$: Eigenvalues of $(\mathbf{R}')^{-1}\mathbf{K}(\mathbf{R}')^{-1}$
 \mathbf{Q} : Orthogonal matrix
- (iii.) Then set transformed design matrix $\tilde{\mathbf{Z}} = \mathbf{Z}(\mathbf{R}')^{-1}\mathbf{Q}$ such that
 $\mathbf{f} = \mathbf{Z}\boldsymbol{\beta} = \tilde{\mathbf{Z}}\tilde{\boldsymbol{\beta}}$ ($\boldsymbol{\beta} = (\mathbf{R}')^{-1}\mathbf{Q}\tilde{\boldsymbol{\beta}}$)
- (iv.) and the resulting penalty is now given by
 $\boldsymbol{\beta}'\mathbf{K}\boldsymbol{\beta} = \tilde{\boldsymbol{\beta}}'\mathbf{Q}'(\mathbf{R}')^{-1}\mathbf{K}(\mathbf{R}')^{-1}\mathbf{Q}\tilde{\boldsymbol{\beta}} = \tilde{\boldsymbol{\beta}}'\mathbf{S}\tilde{\boldsymbol{\beta}}$

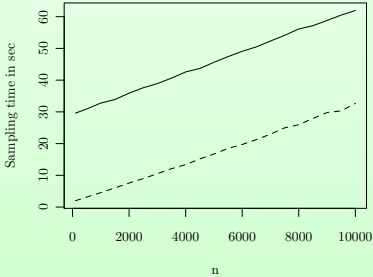
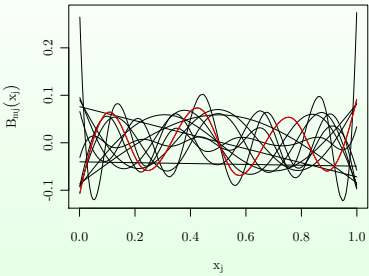
Mean and precision matrix are now given by

$$\mu_{\tilde{\beta}_{mj}} = \frac{1}{1 + \lambda_j s_{mj}} \cdot u_{mj} \quad m = 1, \dots, M_j$$

where $\lambda_j = \sigma^2 / \tau_j^2$ and u_{mj} is the m -th element of the vector $\mathbf{u}_j = \tilde{\mathbf{Z}}_j \mathbf{W} (\mathbf{y} - \boldsymbol{\eta} + \mathbf{f}_j)$, and entries of the corresponding diagonal precision matrix

$$\Sigma_{\tilde{\beta}_j}[m, m] = \frac{\sigma^2}{1 + \lambda_j s_{mj}} \quad m = 1, \dots, M_j$$

Alternative sampling scheme based on transformed parametrization



MCMC sampling scheme

for $t = 1, \dots, T$ {

1. for $j = 1, \dots, p$ {

$$1.1 \quad \tilde{\boldsymbol{\beta}}_j^{(t+1)} | \cdot \sim N \left(\boldsymbol{\mu}_{\tilde{\boldsymbol{\beta}}_j}^{(t)}, \boldsymbol{\Sigma}_{\tilde{\boldsymbol{\beta}}_j}^{(t)} \right)$$

1.2 if level within $\tilde{\boldsymbol{\beta}}_j$ set $\mathbf{y}^* = \tilde{\boldsymbol{\beta}}_j^{(t+1)}$ and repeat steps 1-4

$$1.3 \quad \tau_j^{2(t+1)} | \cdot \sim IG \left(a + \frac{rk(\mathbf{K}_j)}{2}, b + \frac{1}{2} \tilde{\boldsymbol{\beta}}_j'^{(t+1)} \mathbf{K}_j \tilde{\boldsymbol{\beta}}_j^{(t+1)} \right)$$

1.4 update $\boldsymbol{\eta}$

}

$$2. \quad \tilde{\boldsymbol{\gamma}}^{(t+1)} | \cdot \sim N \left(\boldsymbol{\mu}_{\tilde{\boldsymbol{\gamma}}}^{(t)}, \boldsymbol{\Sigma}_{\tilde{\boldsymbol{\gamma}}}^{(t)} \right)$$

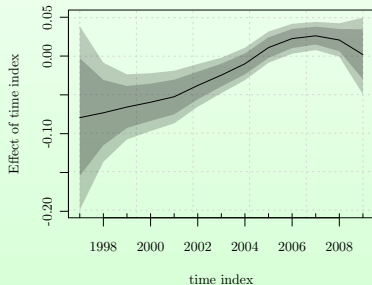
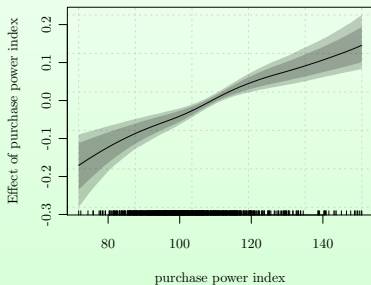
3. update $\boldsymbol{\eta}$

$$4. \quad \sigma^{2(t+1)} | \cdot \sim IG \left(a + \frac{n}{2}, b + \frac{1}{2} (\mathbf{y} - \boldsymbol{\eta}^{(t+1)})' (\mathbf{y} - \boldsymbol{\eta}^{(t+1)}) \right)$$

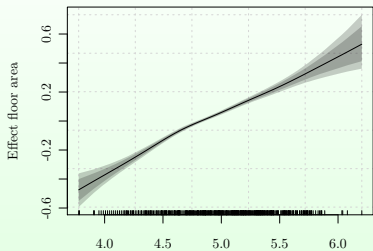
}

Results: Hedonic regression data for house prices

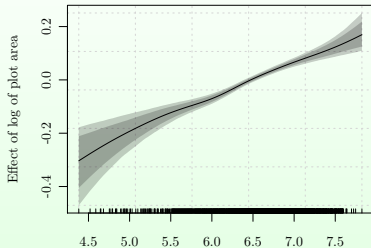
Structural continuous covariates



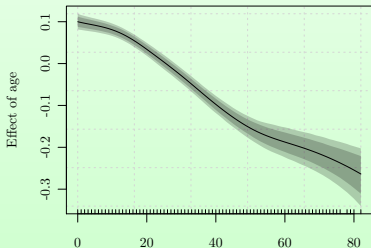
Structural continuous covariates



floor area

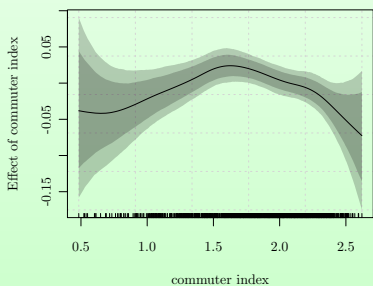
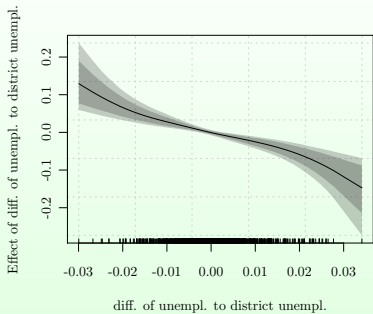
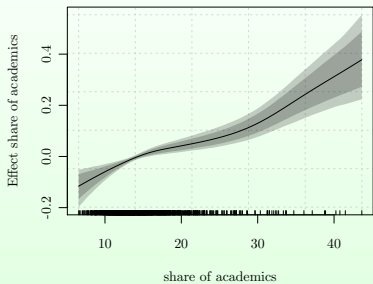


log of plot area

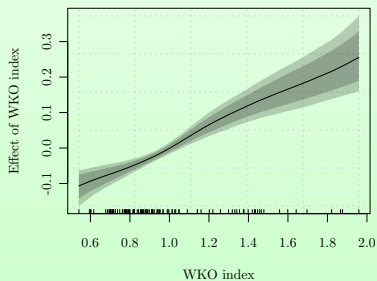
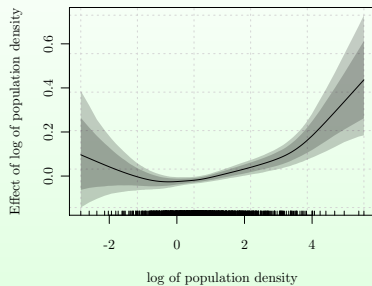
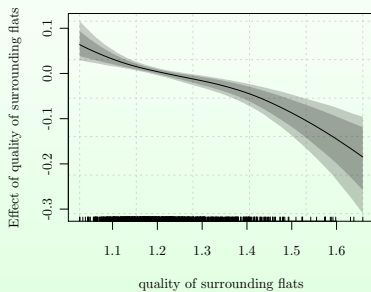


age

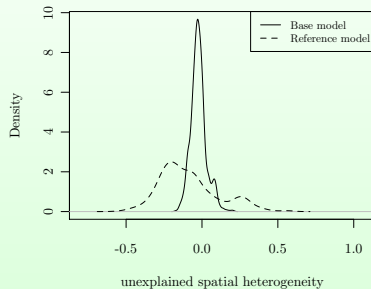
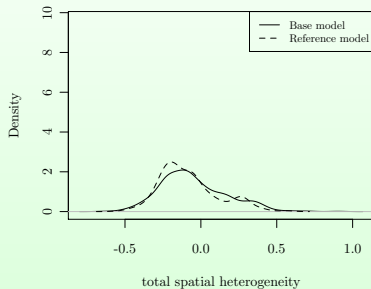
Neighborhood effects



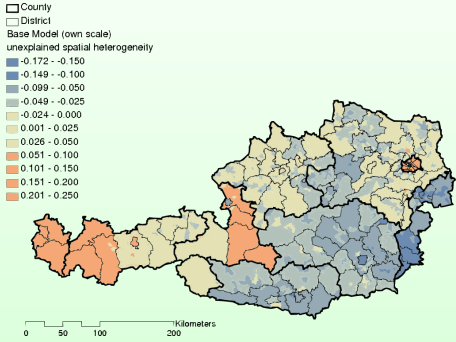
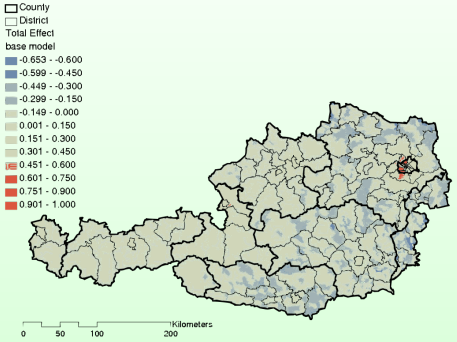
Neighborhood effects



Neighborhood effects



Results: Hedonic regression data for house prices



Thank you!!!

- Belitz, C., Lang, S. (2008). Simultaneous selection of variables and smoothing parameters in structured additive regression models. *Computational Statistics and Data Analysis*, **53**, 61-81.
- Brunauer, W., Lang, S. and Umlauf, N. (2010). Modeling House Prices using Multilevel Structured Additive Regression. Technical report, University of Innsbruck.
- Eilers, P. H. C. and Marx, B. D. (1996). Flexible smoothing using B-splines and penalized likelihood. *Statistical Science*, **11**, 89-121.
- Fahrmeir, L., Kneib, T. and Lang, S. (2004). Penalized structured additive regression for space-time data: a Bayesian perspective. *Statistica Sinica*, **14**, 731-761.
- Lang, S., Umlauf, N., Kneib, T., Hartgen, K. and Wechselberger, P. (2010): Multilevel Generalized Structured Additive Regression. Technical report, University of Innsbruck