

Models for autocorrelated data

Vojtěch Barták

LS 2022

- 1 General overview
- 2 Models with correlated errors
- 3 Models with uncorrelated errors

Section 1

General overview

General model for spatial data

$$\mathbf{Y}(s) = \mu(s) + \mathbf{e}(s)$$

\mathbf{Y} ... vector of response observations

s ... vector of spatial coordinates

μ ... deterministic mean function

\mathbf{e} ... random “error” component

The spatial structure observed in \mathbf{Y} can be modelled:

- in the **mean component**
 - Spatially structured predictors (induced spatial structure)
 - Autocovariate models
 - Trend surface models
 - Moran’s eigenvector mapping
- in the **random component**
 - Geostatistical models
 - Autoregressive models
- both

Modeling spatial structure in the mean component

Linear models:

$$\boldsymbol{\mu}(\mathbf{s}) = \mathbf{X}(\mathbf{s})\boldsymbol{\beta}, \mathbf{e}(\mathbf{s}) \sim MVN(\mathbf{0}, \sigma^2 \mathbf{I})$$

$\mathbf{X}(\mathbf{s})$... matrix of fixed predictors, including spatially structured ones

$\boldsymbol{\beta}$... vector of unknown parameters (fixed effects)

- estimation by *ordinary least squares* (OLS)

Generalized linear models:

- 1 Y_i are mutually independent, following a common distributional family (Gaussian, Poisson, Binomial, ...)
- 2 $\boldsymbol{\mu}(\mathbf{s}) = u[\mathbf{X}(\mathbf{s})\boldsymbol{\beta}]$, u ... link function

- estimation by maximum likelihood (ML)

Modeling spatial structure in the random component

$$\mathbf{e} \sim (\mathbf{0}, \mathbf{\Sigma}(\boldsymbol{\theta}))$$

- $\mathbf{\Sigma}$ is a positive definite matrix with at least some non-zero off-diagonal elements
- $\boldsymbol{\theta}$ is a vector of parameters describing the spatial dependence
- Trying to capture the nature of spatial dependence - the real spatial autocorrelation
- Relies on **stationarity** assumption
- Fixed effects can be estimated by *generalized least squares* (GLS)
- Can be viewed as **mixed models** (estimation by ML/REML)

Section 2

Models with correlated errors

Geostatistical linear model

$$\mu(\mathbf{s}) = \mathbf{X}(\mathbf{s})\beta$$

$$\mathbf{e}(\mathbf{s}) = \mathbf{S}(\mathbf{s}) + \epsilon(\mathbf{s})$$

$\mathbf{S}(\mathbf{s}) \sim MVN(\mathbf{0}, \mathbf{C}) \dots$ Gaussian process

$\mathbf{C} = (\sigma_{i,j}), \sigma_{i,j} = \text{Cov}(S_i, S_j) = C(u) \dots$ Covariance function

$\epsilon(\mathbf{s}) \sim MVN(\mathbf{0}, \tau^2 \mathbf{I}) \dots$ Nugget effect

$$\mathbf{\Sigma} = \mathbf{C} + \tau^2 \mathbf{I}$$

Covariance function estimated from data:

- by fitting the curve to the sample variogram (*classical geostatistics*)
- by ML/REML techniques together with other parameters
(*model-based approach*)

Geostatistical linear model

Spatial prediction/interpolation: **kriging**

$$\hat{Y}(\mathbf{s}_0) = \mathbf{x}'\hat{\beta}_{GLS} + \mathbf{c}'\Sigma^{-1}(\mathbf{Y} - \mathbf{X}\hat{\beta}_{GLS})$$

$\mathbf{c} = (C(\mathbf{s}_0, \mathbf{s}_1), C(\mathbf{s}_0, \mathbf{s}_2), \dots, C(\mathbf{s}_0, \mathbf{s}_n))'$

$\mathbf{c}'\Sigma^{-1} \dots$ *kriging weights*

$\text{var}(\hat{Y}(\mathbf{s}_0)) \dots$ *kriging variance*

- Simple kriging ... known constant mean
- Ordinary kriging ... unknown constant mean
- Universal kriging ... unknown mean depending on covariates

Geostatistical GLM

Generalized linear model (GLM):

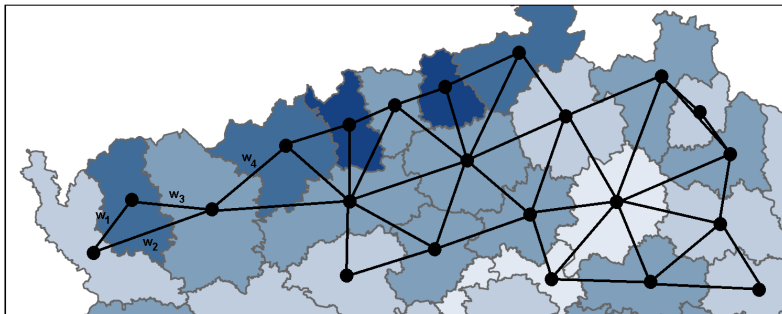
- ① $Y_i \dots$ mutually independent Gaussian/Poisson/Binomial variables
- ② $\mu(\mathbf{s}) = u[\mathbf{X}(\mathbf{s})\beta]$, $u \dots$ link function

Generalized linear geostatistical model (GLGM):

- ① $Y_i \dots$ mutually independent Gaussian/Poisson/Binomial variables
 - ② $\mu(\mathbf{s}) = u[\mathbf{X}(\mathbf{s})\beta + \mathbf{S}(\mathbf{s})]$, $u \dots$ link function
 - ③ $\mathbf{S}(\mathbf{s}) \dots$ Gaussian process with zero mean and some covariance function
- Special case of GLMM

Autoregressive models

- Based on discrete locations with a **neighborhood structure**
- Magnitude of spatial interactions between neighbors -> **spatial weights**



- Not necessarily areal data...

Simultaneous autoregressive model

$$\mathbf{Y}(\mathbf{s}) = \mathbf{X}(\mathbf{s})\boldsymbol{\beta} + \mathbf{e}(\mathbf{s})$$

$$\mathbf{e}(\mathbf{s}) = \mathbf{B}\mathbf{e}(\mathbf{s}) + \boldsymbol{\epsilon}(\mathbf{s})$$

\mathbf{B} ... matrix of spatial dependence parameters, $b_{i,i} = 0$

$$\boldsymbol{\epsilon}(\mathbf{s}) \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$\boldsymbol{\Sigma}_{SAR} = (\mathbf{I} - \mathbf{B})^{-1} \sigma^2 \mathbf{I} (\mathbf{I} - \mathbf{B}')^{-1}$$

Usually:

$$\mathbf{B} = \rho \mathbf{W}$$

ρ ... single correlation parameter

\mathbf{W} ... matrix of spatial weights

Conditional autoregressive model

Assumption: Spatial process is *Markov random field*

$$E[Y(\mathbf{s}_i) | \mathbf{Y}(\mathbf{s})_{-i}] = \mathbf{x}(\mathbf{s}_i)' \boldsymbol{\beta} + \sum_{j=1}^n b_{i,j} [Y(\mathbf{s}_j) - \mathbf{x}(\mathbf{s}_j)' \boldsymbol{\beta}]$$

$$\text{Var}[Y(\mathbf{s}_i) | \mathbf{Y}(\mathbf{s})_{-i}] = \sigma^2$$

- $\boldsymbol{\Sigma}_{CAR} = (\mathbf{I} - \mathbf{B})^{-1} \sigma^2$
- Again, usually $\mathbf{B} = \rho \mathbf{W}$
- In these constant-variance versions, every SAR model can be expressed as a CAR model, and vice versa

Section 3

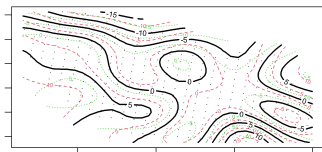
Models with uncorrelated errors

Autocovariate models

- Each observation is modeled as depending on a summary of neighboring observations
- This *autocovariate* is prepared first and added as a fixed predictor
- A “naive” approach to autoregression
- Typically a weighted average of the neighboring values
- Can be used in any type of model (G)L(M)M, GAM, Random Forests, . . .)
- Assumes stationarity

Trend surface models

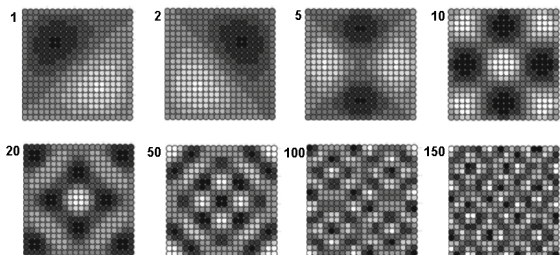
- Originally: a simple (linear, quadratic) spatial trend is added as fixed predictor
- Combined with geostatistical models
- Extended to complex, nonlinear smooth surfaces, describing the unexplained spatial structure
- Typically a smooth term $s(x, y)$ in a **generalized additive model (GAM)**



- Description, but not explanation!
- **Does not assume stationarity!**

Moran's eigenvectors mapping

- PCA applied to the distance or weight matrix
- The resulting variables (called Moran's eigenvectors) used as fixed predictors
- Only those associated with positive eigenvalues (positive spatial dependence)



- Similar to trend surface models and Fourier transform
- Significant eigenvectors represent spatial dependence at different scales