

# Lecture 5 Hierarchical modeling of spatial data

Shiwei Lan<sup>1</sup>

<sup>1</sup>School of Mathematical and Statistical Sciences  
Arizona State University

STP598 Spatiotemporal Analysis  
Fall 2020

Hierarchical  
Model

S.Lan

Stationary  
spatial process  
models

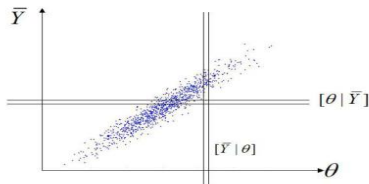
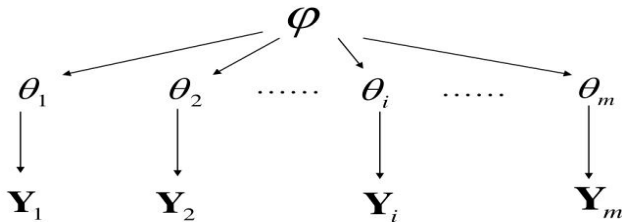
Generalized  
linear spatial  
process  
modeling

Areal data  
models

General linear  
areal data  
modeling

- 1 Stationary spatial process models
- 2 Generalized linear spatial process modeling
- 3 Areal data models
- 4 General linear areal data modeling

## Hierarchical model



- The basic model is

$$Y(s) = \mu(s) + w(s) + \epsilon(s) \quad (1)$$

- The mean structure  $\mu(s) = x^T \beta$ .
- The residual is partitioned into two parts: the spatial part  $w(s)$  and the non-spatial errors  $\epsilon(s)$ .
- Recall that the  $w(s)$  introduces the partial sill ( $\sigma^2$ ) and the range  $\phi$ ; while the  $\epsilon(s)$  adds the nugget  $\tau^2$ .
- Assume stationarity. Pure error  $\epsilon(s)$  vs spatial error  $w(s)$ . We further assume
  - $w(s+h) - w(s) \rightarrow 0$  as  $h \rightarrow 0$
  - $[w(s+h) + \epsilon(s+h)] - [w(s) + \epsilon(s)] \not\rightarrow 0$  as  $h \rightarrow 0$
- *Microscale* view:  $\epsilon(s)$  is a spatial process with very rapid decay in association, and only matters at high resolution.

- Suppose we have data  $Y(s_i)$ ,  $i = 1, \dots, n$ , and let  $Y = (Y(s_1), \dots, Y(s_n))^T$ . In the basic Gaussian isotropic kriging model, we assume

$$\Sigma = \sigma^2 H(\phi) + \tau^2 I \quad (2)$$

where  $H$  is correlation matrix with  $H_{ij} = \rho(s_i - s_j; \phi)$  and  $\rho$  is a valid isotropic correlation function on  $\mathbb{R}^2$ , e.g.  $\rho(s_i - s_j; \phi) = \exp(-\phi \|s_i - s_j\|)$

- Collecting all model parameters into a vector  $\theta = (\beta, \sigma^2, \tau^2, \phi)$ , we have

$$p(\theta|y) \propto f(y|\theta)p(\theta) \quad (3)$$

where

$$Y|\theta \sim N(X\beta, \sigma^2 H(\phi) + \tau^2 I) \quad (4)$$

- Typically, independent priors are chosen for the different parameters

$$p(\boldsymbol{\theta}) = p(\boldsymbol{\beta})p(\sigma^2)p(\tau^2)p(\phi) \quad (5)$$

- We usually adopt the following priors

$$\boldsymbol{\beta} \sim N_{p+1}(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0) \quad (6)$$

$$\sigma^2 \sim \Gamma^{-1}(a_1, b_1) \quad (7)$$

$$\tau^2 \sim \Gamma^{-1}(a_2, b_2) \quad (8)$$

$$\phi \sim \Gamma(a_3, b_3) \quad (9)$$

- In Matérn class, the product  $\sigma^2\phi^{2\nu}$  can be identified but not the individuals (Zhang, 2004). Often, a very informative prior is imposed for  $\phi$  (e.g. uniform over interval) and a relatively vague prior is used for  $\sigma^2$ .

- When  $\tau^2 \equiv 0$ , with the exponential covariance, Berger et al. (2001) impose the *objective* priors of the form

$$p(\boldsymbol{\beta}, \sigma^2, \phi) \propto \frac{p(\phi)}{(\sigma^2)^\alpha} \quad (10)$$

- This implies a flat prior for  $\boldsymbol{\beta}$ . With a uniform prior for  $\phi$ , it can be shown that an improper posterior arises for  $\alpha < 2$ .
- If  $\sigma^2 \sim \Gamma^{-1}(\epsilon, \epsilon)$ , then  $\alpha = 1 + \epsilon$  getting an improper posterior for small  $\epsilon$ .
- If  $\sigma^2 \sim \Gamma^{-1}(a_1, b_1)$ ,  $a_1 \geq 1$  is recommended  $\alpha \geq 2$ .
- Closed form of *marginal* posterior of  $\boldsymbol{\beta}$  may be available (e.g. with Gaussian likelihood).

- The hierarchical modeling shares the following generic format

$$[\text{data}|\text{process, parameters}] [\text{process}|\text{parameters}] [\text{parameters}] \quad (11)$$

- The previous spatial process model (4) can be recast as a hierarchical model with two stages

$$Y|\theta, W \sim N(X\beta + W, \tau^2 I) \quad (12)$$

$$W|\sigma^2, \phi \sim N(0, \sigma^2 H(\phi)) \quad (13)$$

where  $W = (w(s_1), \dots, w(s_n))^T$  is the spatial random effect that captures spatial dependence.

- The parameter space is now augmented from  $\theta$  to  $(\theta, W)$ , with the dimension increased by  $n$ .



- The resulting  $p(\boldsymbol{\theta}|y)$  is the same, but we have the choice of using MCMC to fit either  $f(y|\boldsymbol{\theta})p(\boldsymbol{\theta})$  or  $f(y|\boldsymbol{\theta}, W)p(W|\boldsymbol{\theta})p(\boldsymbol{\theta})$ .
- Interest is often in the spatial surface  $W|y$  as well as prediction for  $W(s_0)|y$  for various choices of  $s_0$ .
- Note  $p(W|y)$  can be recovered from  $p(\boldsymbol{\theta}|y)$ :

$$p(W|y) = \int p(W|\boldsymbol{\theta}, y)p(\boldsymbol{\theta}|y)d\boldsymbol{\theta} \quad (14)$$

which can be sampled by one for one composition of posterior sampling:  
 $W^{(g)} \sim p(W|\boldsymbol{\theta}^{(g)}, y)$  with  $\boldsymbol{\theta}^{(g)} \sim p(\boldsymbol{\theta}|y)$ .

- Bayesian kriging involves prediction of  $Y_0 \equiv Y(s_0)$  at a new location  $s_0$  with associated covariate  $x(s_0)$ :

$$p(y_0|y, X, x_0) = \int p(y_0, \boldsymbol{\theta}|y, X, x_0)d\boldsymbol{\theta} = \int p(y_0|y, \boldsymbol{\theta}, x_0)p(\boldsymbol{\theta}|y, X)d\boldsymbol{\theta} \quad (15)$$

- In practice, we use MCMC to obtain posterior samples  $\{\boldsymbol{\theta}^{(g)}\}_{g=1}^G$  with  $\boldsymbol{\theta}^{(g)} \sim p(\boldsymbol{\theta}|y, X)$ , and approximate the above predictive distribution with

$$\hat{p}(y_0|y, X, x_0) = \frac{1}{G} \sum_{g=1}^G p(y_0|y, \boldsymbol{\theta}^{(g)}, x_0) \quad (16)$$

- Replicates of prediction  $y_0^{(g)}$  can be obtained using the composition sampling.
- Multi-output prediction for  $Y_0 = (Y(s_{01}), \dots, Y(s_{0m}))^T$  is also available

$$p(y_0|y, X, x_0) = \int p(y_0|y, \boldsymbol{\theta}, x_0)p(\boldsymbol{\theta}|y, X)d\boldsymbol{\theta} \approx \frac{1}{G} \sum_{g=1}^G p(y_0|y, \boldsymbol{\theta}^{(g)}, x_0) \quad (17)$$

Hierarchical  
Model

S.Lan

Stationary  
spatial process  
models

Generalized  
linear spatial  
process  
modeling

Areal data  
models

General linear  
areal data  
modeling

- 1 Stationary spatial process models
- 2 Generalized linear spatial process modeling
- 3 Areal data models
- 4 General linear areal data modeling

- In some point-referenced data sets, measurements  $Y(s)$  are not naturally modeled as a normal distribution. e.g. binary data, and counting data.
- The observations  $Y(s_i)$  are modeled independent conditioned on  $\beta$  and  $w(s_i)$  within the class of exponential family:

$$p(y(s_i)|\beta, w(s_i), \gamma) = h(y(s_i), \gamma) \exp \{ \gamma [y(s_i)\eta(s_i) - \psi(\eta(s_i))] \} \quad (18)$$

where  $g(\eta(s_i)) = x^T(s_i)\beta + w(s_i)$  for some link function  $g$ , and  $\gamma$  is a dispersion parameter.

- We presume the  $w(s_i)$  to be the spatial random effect coming from a Gaussian process, i.e.  $W \sim N(0, \sigma^2 H(\phi))$ .

- Using conditional independence, we have the joint distribution

$$f(y(s_1), \dots, y(s_n) | \beta, \sigma^2, \phi, \gamma) = \int \left( \prod_{i=1}^n f(y(s_i) | \beta, w(s_i), \gamma) \right) p(W | \sigma^2, \phi) dW \quad (19)$$

- Note  $W$  may not be marginalized out in general.
- Binary data.** We model  $Y(s)$  through the latent process  $Z(s) = x(s)^T \beta + w(s) + \epsilon(s)$ :

$$\Pr(Y(s) = 1) = \Pr(Z(s) \geq 0) = g^{-1}(x(s)^T \beta + w(s)) \quad (20)$$

for some link function  $g(\cdot)$  such that  $g^{-1}$  takes  $p(s) := \Pr(Y(s) = 1)$  to  $\mathbb{R}^1$ .

- Possible choices: the logit  $g(x) = \log \frac{x}{1-x}$  and the probit  $g(x) = \Phi^{-1}(x)$ .

# Generalized linear spatial process modeling

Hierarchical  
Model

S.Lan

Stationary  
spatial process  
models

Generalized  
linear spatial  
process  
modeling

Areal data  
models

General linear  
areal data  
modeling

- Consider real estate data at 50 locations in Baton Rouge, LA.
- $Y(s) = 1$  indicates that the price of the property at location  $s$  is “high” (above the median for the region);  $Y(s) = 0$  indicates that the price is “low”.
- Covariates: the house’s age, total living area, and other area in the property.
- We fit the model with the logit link, assuming vague priors for  $\beta$ , a  $\text{unif}(0, 10)$  prior for  $\phi$  and a  $\Gamma^{-1}(0.1, 0.1)$  prior for  $\sigma^2$ .

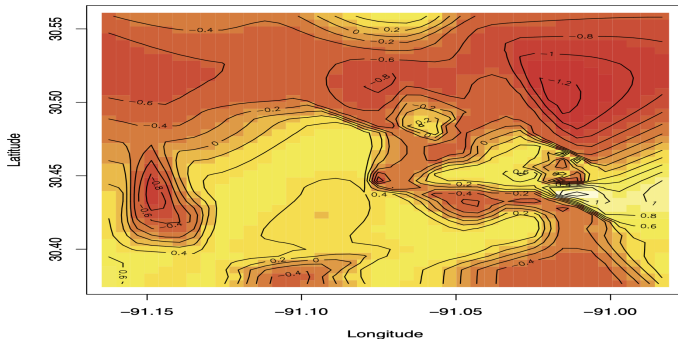


Figure 6.7 Image plot of the posterior median surface of the latent spatial process  $w(s)$ , binary 14 / 27

- **Counting data.** We model  $Y(s)$  through the latent process  $Z(s) = x(s)^T \beta + w(s) + \epsilon(s)$ :

$$Y(s) \sim \text{pois}(g^{-1}(Z(s))) \quad (21)$$

for some link function  $g(\cdot)$ , e.g. the canonical link  $g(x) = \log(x)$ .

- This is related to the log-Gaussian Cox Process (LGCP) model.

Hierarchical  
Model

S.Lan

Stationary  
spatial process  
models

Generalized  
linear spatial  
process  
modeling

Areal data  
models

General linear  
areal data  
modeling

- 1 Stationary spatial process models
- 2 Generalized linear spatial process modeling
- 3 Areal data models
- 4 General linear areal data modeling



- An area of strong biostatistical and epidemiological interest is that of *disease mapping*.
- We typically have count data of the following:

$Y_i$  = observed number of cases of disease in county  $i, i = 1, \dots, I$

$E_i$  = expected number of cases of disease in county  $i, i = 1, \dots, I$

where  $E_i$  are thought of as fixed and known functions of  $n_i$ , the number of persons at risk for the disease in county  $i$ .

- One can simply assume

$$E_i = n_i \bar{r} = n_i \frac{\sum_i y_i}{\sum_i n_i} = \sum_i y_i \frac{n_i}{\sum_i n_i} \quad (22)$$

i.e.  $\bar{r}$  is the overall disease rate in the entire study region.

- However, such *internal standardization* is “cheating”. And one might modify it by age-adjusted rates for the disease

$$E_i = n_{ij}r_j \quad (23)$$

where  $n_{ij}$  is the person-years at risk in area  $i$  for age group  $j$ , and  $\bar{r}_j$  is the disease rate in age group  $j$ . This process is called *external standardization*.

- The usual model for  $Y_i$  is the Poisson model

$$Y_i | \eta_i \sim \text{Pois}(E_i \eta_i) \quad (24)$$

where  $\eta_i$  is the true *relative risk* of disease in region  $i$ .

- The maximum likelihood estimate (MLE) of  $\eta_i$  is

$$\hat{\eta}_i = \text{SMR}_i = \frac{Y_i}{E_i} \quad (25)$$

- $\text{Var}(\text{SMR}_i) = \text{Var}(Y_i)/E_i^2 = \eta_i/E_i$ .  $\widehat{\text{Var}}(\text{SMR}_i) = \hat{\eta}_i/E_i = Y_i/E_i^2$ .

- For detecting extra-Poisson variability (overdispersion) in the observed rates, we seek *random effects* model for  $\eta_i$  through hierarchical Bayesian modeling.
- Poisson-gamma model

$$Y_i | \eta_i \stackrel{ind}{\sim} \text{pois}(E_i \eta_i), \quad i = 1, \dots, I \quad (26)$$

$$\eta_i \stackrel{iid}{\sim} \Gamma(a, b) \quad (27)$$

- Choose  $a, b$  based on  $\mu = a/b$  and  $\sigma^2 = a/b^2$ .
- Due to conjugacy, we have  $p(\eta_i | y_i) = \Gamma(a + y_i, b + E_i)$ .
- The Bayesian posterior mean (point estimate) of  $\eta_i$  is

$$E(\eta_i | y) = E[\eta_i | y_i] = \frac{y_i + a}{E_i + b} = \frac{y_i + \frac{\mu^2}{\sigma^2}}{E_i + \frac{\mu}{\sigma^2}} = w_i \text{SMR}_i + (1 - w_i) \mu \quad (28)$$

$$\text{where } w_i = \frac{E_i}{E_i + \frac{\mu}{\sigma^2}}.$$

- The gamma prior suffers from a serious defect: it fails to allow for spatial correlation among  $\eta_i$ 's. We introduce the following model to overcome this issue.
- **Poisson-lognormal model.** Denote  $\psi_i \equiv \log \eta_i$ , the log-relative risks.

$$Y_i | \psi_i \stackrel{ind}{\sim} \text{pois}(E_i e^{\psi_i}), \quad \psi_i = \mathbf{x}_i' \boldsymbol{\beta} + \theta_i + \phi_i \quad (29)$$

$$\theta_i \stackrel{iid}{\sim} N(0, 1/\tau_h) \quad (30)$$

$$\boldsymbol{\phi} | \boldsymbol{\mu}, \boldsymbol{\lambda} \sim N_I(\boldsymbol{\mu}, H(\boldsymbol{\lambda})), \quad \boldsymbol{\phi} = (\phi_1, \dots, \phi_I)' \quad (31)$$

- $\theta_i$ 's capture region-wide *heterogeneity*. These random effects capture extra-Poisson variability in the log-relative risks that varies “globally”.
- $\boldsymbol{\phi}$  capture regional *clustering*. They model extra-Poisson variability in the log-relative risks that varies “locally”.

- In specifying the covariance  $H(\lambda)$ , what are appropriate inter-areal unit distances? Should we use centroid to centroid? Does this make sense with units of quite differing sizes and irregular shapes?
- Big  $N$  problem emerges with regard to high-dimensional matrix inversion. We return to CAR (IAR) specifications for  $\phi$ :

$$\phi \sim \text{CAR}(\tau_c), \quad p(\phi) \propto \exp \left\{ -\frac{\tau_c}{2} \sum_{i \neq j} w_{ij} (\phi_i - \phi_j)^2 \right\} \quad (32)$$

where  $w_{ij}$  is the 0-1 adjacency weights.

- Full conditional of  $\phi_i$

$$p(\phi_i | \phi_{j \neq i}, \theta, \beta, y) \propto \text{pois}(y_i | E_i e^{x_i' \beta + \theta_i + \phi_i}) N(\phi_i | \bar{\phi}_i, 1/(\tau_c m_i)) \quad (33)$$

- There are two issues: i) impropriety; ii) selection of  $\tau_h$  and  $\tau_c$ .
- Recall that IAR prior is improper due to the singular precision matrix  $\Sigma^{-1} = D_w - W$ . We could consider either  $\Sigma^{-1} = D_w - \rho W$  or  $\Sigma^{-1} = M^{-1}(I - \alpha \widetilde{W})$ .
- Alternatively, one could ignore the impropriety of the standard CAR model and consider the intrinsic CAR. In practice, we could impose the constraint  $\sum_{i=1}^I \phi_i = 0$  and numerically implement this centering on the fly.
- There is an identifiability issue of  $\theta_i$  and  $\phi_i$  in the model. One can choose to fix  $\tau_h$  and  $\tau_c$ .
- Alternatively, one could use gamma priors  $\tau_h \sim \Gamma(a_h, b_h)$  and  $\tau_c \sim \Gamma(a_c, b_c)$  with the following rule of thumb:

$$\text{sd}(\theta_i) = \frac{1}{\sqrt{\tau_h}} \approx \frac{1}{0.7\sqrt{\bar{m}\tau_c}} = \text{sd}(\phi_i) \quad (34)$$

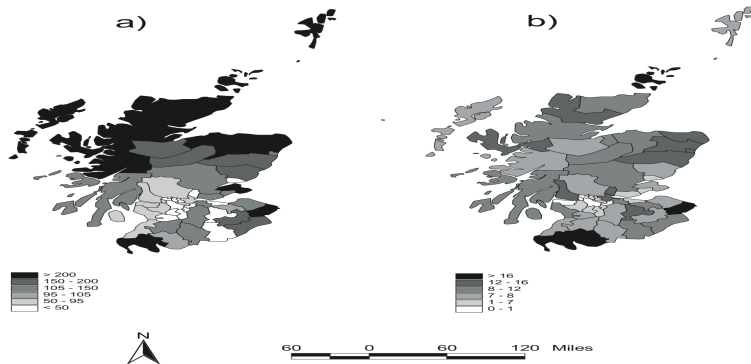


Figure 6.15 *Scotland lip cancer data: (a) crude standardized mortality ratios (observed/expected  $\times 100$ ); (b) AFF covariate values.*

- Consider data from Clayton and Kaldor (1987) are the observed ( $Y_i$ ) and expected ( $E_i$ ) cases of lip cancer for the  $I = 56$  districts of Scotland during the period 1975–1980.
- One county-level covariate  $x_i$ , the percentage of the population engaged in agriculture, fishing or forestry (AFF), is also available.

- We fit the model of log-relative risk

$$\psi_i = \beta_0 + \beta_1 x_i + \theta_i + \phi_i \quad (35)$$

- Note that  $y_i$  cannot inform about  $\theta_i$  or  $\phi_i$ , but only about their sum  $\theta_i + \phi_i$ .
- We are more interested in  $\alpha$ , the proportion of the variability in the random effects that is due to clustering, and choose priors such that  $\alpha \approx 1/2$

$$\alpha = \frac{\text{sd}(\phi)}{\text{sd}(\theta) + \text{sd}(\phi)} \quad (36)$$

Priors for $\tau_c, \tau_h$	Posterior for $\alpha$			Posterior for $\beta$		
	mean	sd	llacf	mean	sd	llacf
G(1.0, 1.0), G(3.2761, 1.81)	.57	.058	.80	.43	.17	.94
G(.1, .1), G(.32761, .181)	.65	.073	.89	.41	.14	.92
G(.1, .1), G(.001, .001)	.82	.10	.98	.38	.13	.91
Priors for $\tau_c, \tau_h$	Posterior for $\xi_1$			Posterior for $\xi_{56}$		
	mean	sd	llacf	mean	sd	llacf
G(1.0, 1.0), G(3.2761, 1.81)	.92	.40	.33	-.96	.52	.12
G(.1, .1), G(.32761, .181)	.89	.36	.28	-.79	.41	.17
G(.1, .1), G(.001, .001)	.90	.34	.31	-.70	.35	.21

Table 6.6 *Posterior summaries for the spatial model with Gamma hyperpriors for  $\tau_c$  and  $\tau_h$ , Scotland lip cancer data; “sd” denotes standard deviation while “llacf” denotes lag 1 sample autocorrelation.*



Hierarchical  
Model

S.Lan

Stationary  
spatial process  
models

Generalized  
linear spatial  
process  
modeling

Areal data  
models

General linear  
areal data  
modeling

- 1 Stationary spatial process models
- 2 Generalized linear spatial process modeling
- 3 Areal data models
- 4 General linear areal data modeling

- Again formulating a hierarchical model,  $Y_i$  may be described using a suitable first-stage member of the exponential family.
- Given  $\beta$  and  $\phi_i$ ,  $Y_i$  is modeled conditionally independent:

$$p(y_i|\beta, \phi_i, \gamma) = h(y_i, \gamma) \exp \{ \gamma [y_i \eta_i - \psi(\eta_i)] \} \quad (37)$$

where  $g(\eta_i) = x_i^T \beta + \phi_i$  for some link function  $g$  with  $\gamma$  is a dispersion parameter.

- The  $\phi_i$ 's will be spatial random effects coming from a CAR model; the pairwise difference, intrinsic (IAR) form is most commonly used.
- In general there is no need to introduce independent heterogeneity effects  $\theta_i$ 's in this generalized linear mixed model with random spatial effects.

- Point-referenced data models are defined with regard to an uncountable number of random variables  $Y(s)$ . For areal units, we envision only a single, finite,  $n$ -dimensional distribution for the  $Y_i$ ,  $i = 1, \dots, I$ .
- With point-referenced data  $Y = (Y(s_1), \dots, Y(s_n))'$ , we model association directly by covariance  $\Sigma_Y$ . With areal data  $Y = (Y_1, \dots, Y_n)'$  and CAR (or SAR) specifications, we instead model the precision  $\Sigma_Y^{-1}$  directly.
- $\Sigma_Y$  encodes *unconditional* association structure; while  $\Sigma_Y^{-1}$  provides *conditional* association structure.
- Explanation is a common goal of point-referenced data modeling, but often an even more important goal is spatial prediction or interpolation (i.e., kriging). With areal units, again a goal is explanation, but now often is supplemented by smoothing (MAUP).
- Big  $N$  problem for spatial process model due to  $\Sigma_Y^{-1}$  and  $|\Sigma_Y|$  in the likelihood; not for CAR (full conditionals) or SAR (no matrix inversion in likelihood).