

Computational Techniques for Spatial Logistic Regression

Certain health outcomes depend on environmental factors that vary geographically (e.g. pollutants), and Epidemiological studies accounting for spatial variability are becoming increasingly popular. Often, the outcome measured is binary (e.g. whether or not a certain disease is diagnosed). Logistic regression with a built-in spatial function is useful for assessing risk at an individual level. However, large sample sizes can lead to computational difficulties. Here, we give a synopsis of various computational techniques for spatial logistic regression discussed in Paciorek, 2007¹.

The general model considered is one where the response, Y_i , is a Bernoulli variable with corresponding probability p_i . For example, p_i may represent the probability that individual i is diagnosed with a certain disease. It is assumed that p_i depends on a vector of covariates x_i (e.g. age, sex or other demographic variables), and a location s_i in two-dimensional space. The logit of p_i is represented as a standard linear model of the covariates x_i , plus a function g of the spatial location, parameterized by θ :

$$Y_i \sim \text{Ber}(p(x_i, s_i)),$$

$$\text{logit}(p(x_i, s_i)) = \mathbf{x}_i^T \boldsymbol{\beta} + g(s_i; \boldsymbol{\theta}),$$

Paciorek focuses on various forms and estimation techniques for the smooth spatial function, g .

1.) Generalized Linear Mixed Models (GLMM)

One approach is to treat the spatial function g as a random effect in a GLMM, $\mathbf{g}_s = \mathbf{Z}\mathbf{u}$. If we assume $\text{Cov}(\mathbf{u}) = \sigma_u^2 \mathbf{I}$, then $\text{Cov}(\mathbf{Z}\mathbf{u}) = \sigma_u^2 \mathbf{Z}\mathbf{Z}^t$, so that the matrix \mathbf{Z} defines the spatial covariance structure. It is suggested to construct \mathbf{Z} based on the Matern covariance:

$$C(\tau) = \frac{1}{\Gamma(v)2^{v-1}} \left(\frac{2\sqrt{v}\tau}{\rho} \right)^v \mathcal{K}_v \left(\frac{2\sqrt{v}\tau}{\rho} \right)$$

where τ is the distance between two spatial locations.

To reduce the computational burden, one could build the covariance structure based on a set of pre-specified knot locations κ_k , $k=1, \dots, K$, with $K < n$. \mathbf{Z} is then constructed as $\mathbf{Z} = \boldsymbol{\psi}\boldsymbol{\Omega}^{-1/2}$, where

$$\boldsymbol{\psi} = (C(|\mathbf{s}_i - \mathbf{k}_k|))_{i=1,\dots,n;k=1,\dots,K}, \text{ and}$$

$$\boldsymbol{\Omega} = (C(|\mathbf{k}_i - \mathbf{k}_j|))_{i=1,\dots,K;j=1,\dots,K} .$$

If one fixes ρ and ν in advance, estimation of σ_u^2 is computationally manageable and several IWLS procedures exist for fitting such a GLMM. Kamman and Wand² suggest taking $\nu = 1.5$ and $\rho = \max_{i,j=1,\dots,n} \|s_i - s_j\|$. This simplifies model estimation, but as \mathbf{g}_s is only parameterized by σ_u^2 , fixing ρ and ν sacrifices model flexibility.

2.) Bayesian GLMM

Another approach, similar to (1), also takes a random effect model $\mathbf{g}_s = \mathbf{Z}\mathbf{u}$ but gives a prior distribution $\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \sigma_u^2 \mathbf{I})$ for the basis coefficients. The model can then be fit via MCMC. In practice, one could include ρ and ν of the Matern covariance in the MCMC. However, in this case recalculating the covariance matrix at each iteration can be computationally difficult.

3.) Bayesian Spectral Basis (SB)

Another approach uses basis functions and the spectral representation of a Gaussian Process (GP). Here, the function g is approximated on a grid, $\mathbf{s}^\#$, of size $K=k_1 \times k_2$. Again we represent g as $\mathbf{g}_{s^\#} = \mathbf{Z}\mathbf{u}$, but here we let \mathbf{Z} be a matrix of orthogonal SB functions, while \mathbf{u} is a vector of complex-valued basis coefficients, $u_m = a_m + b_m$. The Fourier basis is suggested for the matrix \mathbf{Z} .

Let the basis coefficients \mathbf{a} and \mathbf{b} have the prior distribution

$$\begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\theta)$$

Then, the diagonal covariance matrix of the coefficients $\boldsymbol{\Sigma}_\theta$ can be constructed based on the spectral density of a GP covariance function. For the Matern covariance, the spectral density, evaluated at frequency $\boldsymbol{\omega}$, is

$$f_{(\rho,\nu)}(\boldsymbol{\omega}) = \frac{\Gamma(\nu + d/2)(4\nu)^\nu}{\pi^{d/2} \Gamma(\nu)(\pi\rho)^{2\nu}} \cdot \left(\frac{4\nu}{(\pi\rho)^2} + \boldsymbol{\omega}^T \boldsymbol{\omega} \right)^{-(\nu+d/2)}$$

For $\mathbf{m}=(m_1, m_2)$, let $V(a_{m_1,m_2}) = V(b_{m_1,m_2}) = \frac{1}{2} f(\omega_{m_1}^1, \omega_{m_2}^2)$.

Each observation location is then mapped to its nearest grid location through a transformation \mathbf{P} , and then multiplied by σ , the standard deviation of the process: $\mathbf{g}_s = \sigma \mathbf{P} \mathbf{g}_{s^\#}$. Model parameters σ , ρ and ν are given prior distributions and estimated via MCMC. The matrix \mathbf{Z} need not be formed at each iteration, and the operation $\mathbf{Z}\mathbf{u}$ can be done quickly using Fast Fourier Transforms, which helps computational efficiency.

Application

The methods discussed are applied to individual-level leukemia and brain cancer data in Kaohsiung, Taiwan. The metropolitan area contains four petroleum/petrochemical complexes, raising concern over possible adverse health effects. Individuals are geocoded to the location where they have resided for the longest time. The leukemia analysis contains 787 individuals, of whom 206 have been diagnosed with leukemia. The brain cancer analysis contains 576 individuals, of which 165 have been diagnosed with brain cancer.

A map of geocoded addresses is shown in Figure 1. Open circles correspond to cases, while closed circles correspond to controls. The boundaries of the four petrochemical complexes are given by a thick black line.

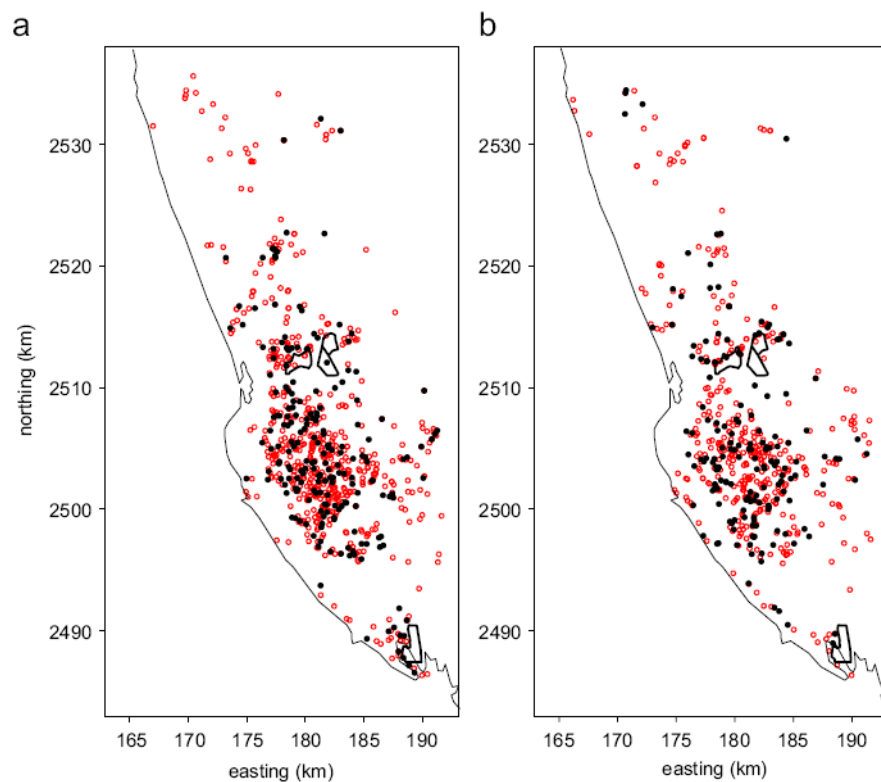


Figure 1: Location of cases and controls for (a) leukemia and (b) brain cancer

Figure 2 gives the posterior mean probability of being a case for each grid cell using the spectral basis method (3). Estimates for both are generally lower in the Northern region, further from the petrochemical complexes, and higher in the South.

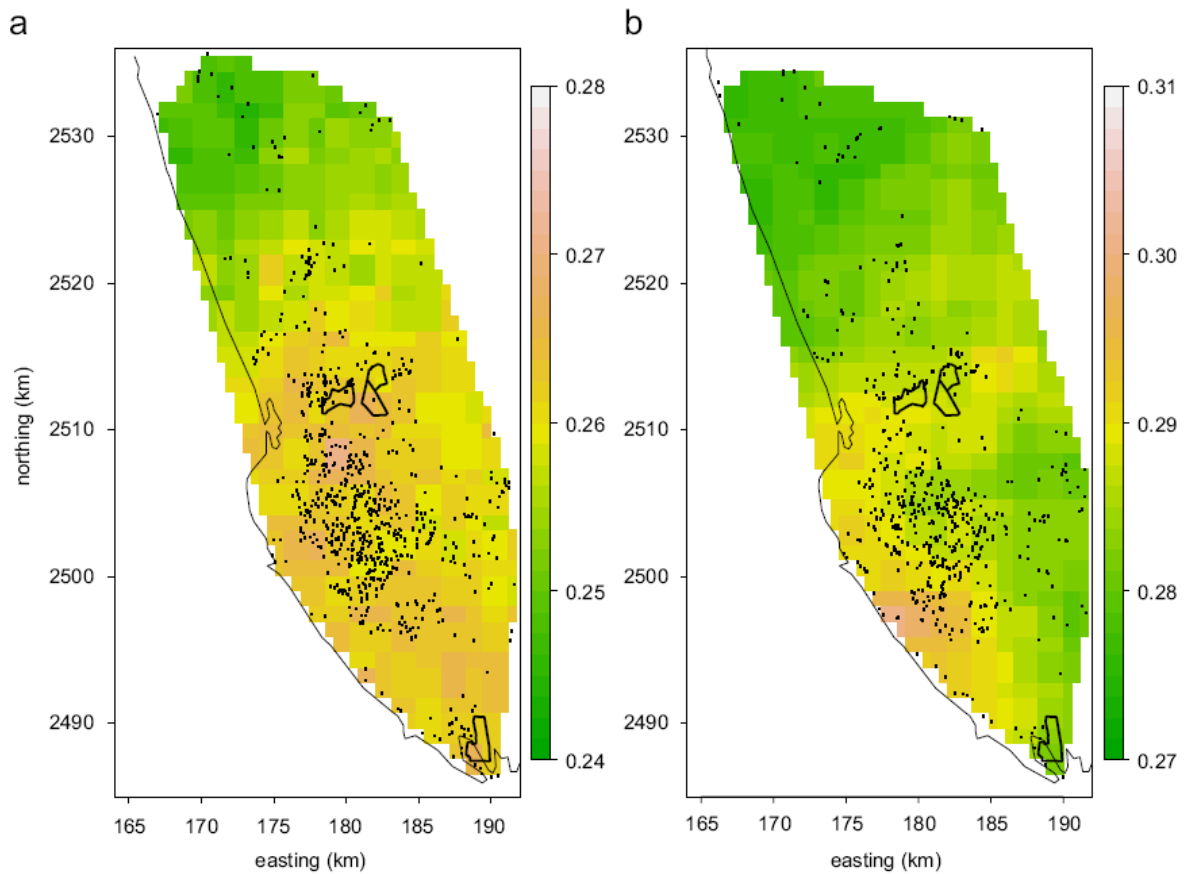


Figure 2: Posterior probability of a case for (a) leukemia and (b) brain cancer

A measure of deviance is used to evaluate various spatial models, including those described in (1) (PL-PQL) and (3) (SB). The model 'Null' corresponds to a non-spatial model. Results are given in Table 1

Table 1: Deviance for various spatial models, and a nonspatial model ('null')

	Leukemia	Brain cancer
PL-GCV	894.5	680.4
PL-PQL	891.7	677.2
SB	891.9	677.9
MRF	892.5	677.8
Null	891.5	678.0

Note that the Null model beats all other models for the leukemia data, and model differences are negligible among the brain cancer data. In fact, the null hypothesis that risk is constant spatially cannot be rejected in either case. This may reflect the absence of a true spatial effect, or an insufficient sample size.

References

¹Paciorek (2007), *Computational Statistics and Data Analysis*, 51, 3631-3653

²Kammann, E., Wand, M. (2003) *Applied Statistics*, 52, 1-18.