

# **Accounting for space in regression models with binary outcomes**

Chris Paciorek

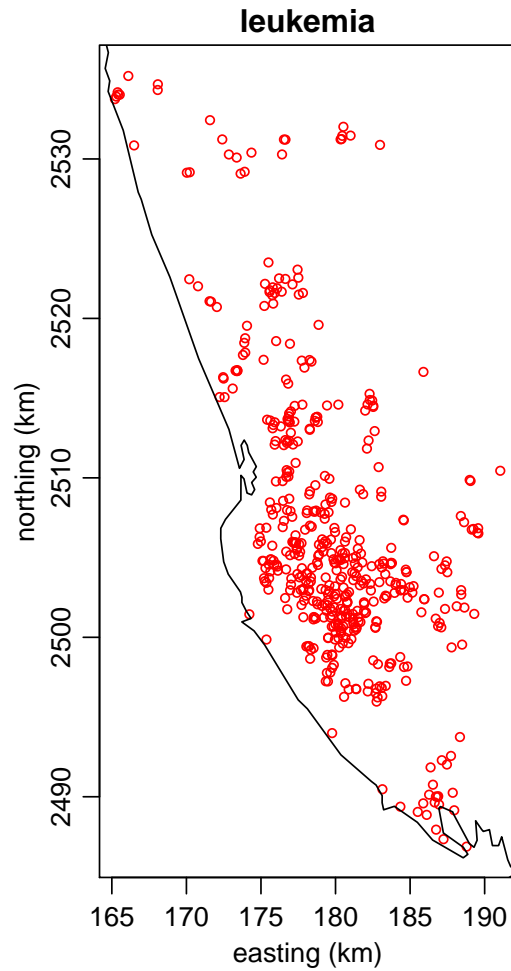
October 15, 2004

Department of Biostatistics

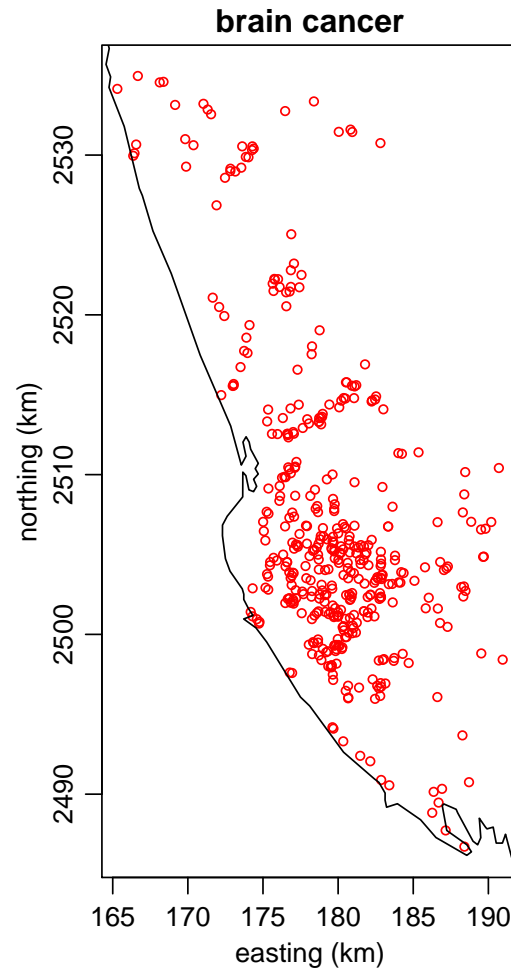
Harvard School of Public Health

[www.biostat.harvard.edu/~paciorek](http://www.biostat.harvard.edu/~paciorek)

# Petrochemical exposure in Kaohsiung, Taiwan

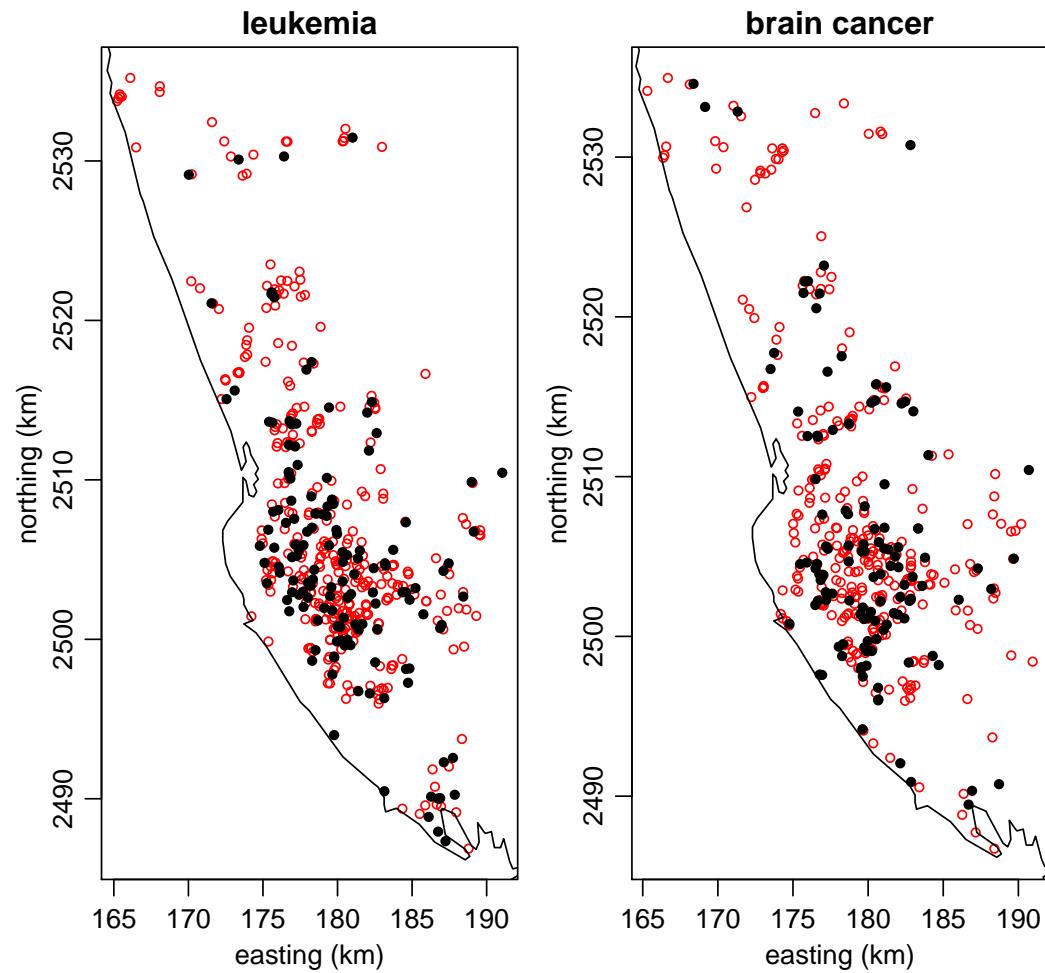


$n = 495$



$n = 433$

# Petrochemical exposure in Kaohsiung, Taiwan



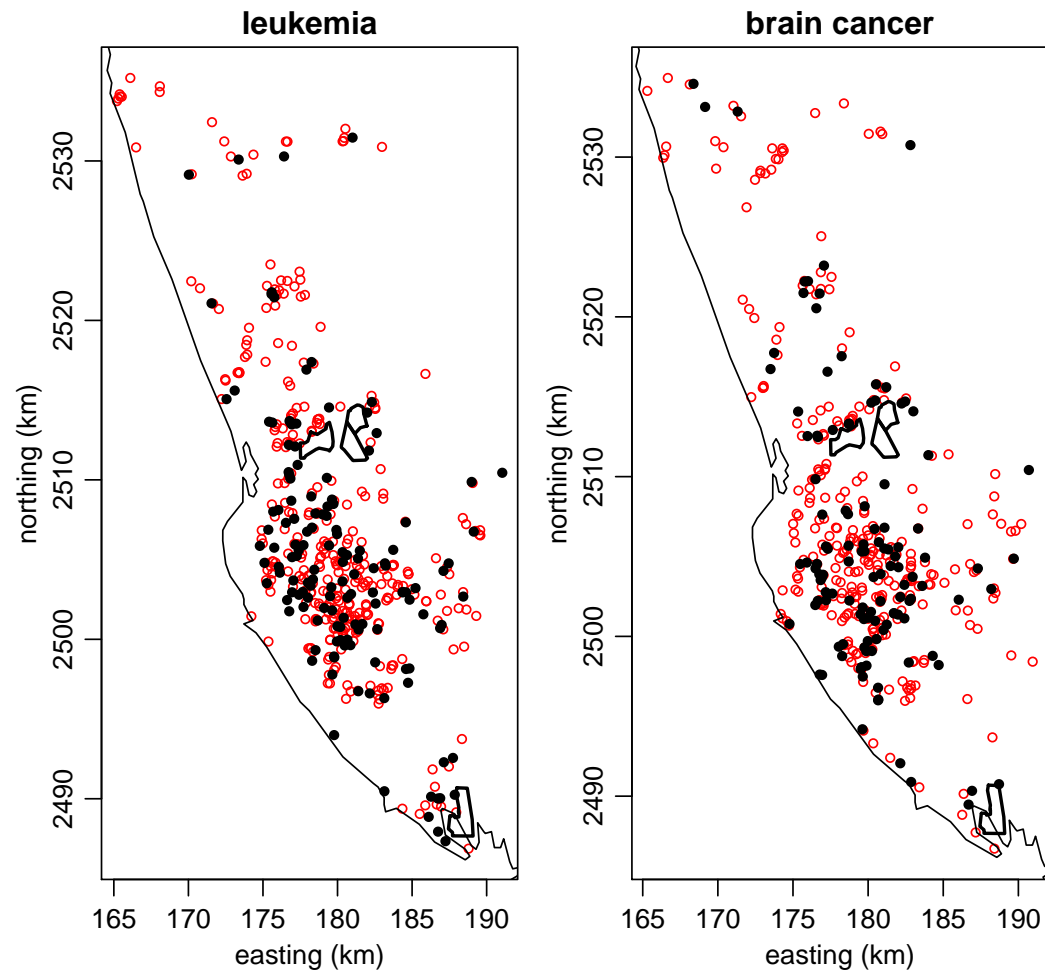
$$n = 495$$

$$n_1 = 141$$

$$n = 433$$

$$n_1 = 121$$

# Petrochemical exposure in Kaohsiung, Taiwan



$$n = 495$$

$$n_1 = 141$$

$$n = 433$$

$$n_1 = 121$$

# Possible approaches for health analysis

- Estimate exposure and use as covariate in health model
- Use distance to exposure source as covariate
- **Explicitly include space as a covariate**
  - **Map of risk - exploratory**
  - **Account for spatially-related unmeasured confounders**
  - **Test for spatial effect**

# Outline

- Motivating example
- Generalized additive model and generalized mixed model approaches
- Difficulties in fitting regression for non-normal outcomes with 2-d smooth terms
- Parameterizations and fitting methods
- Simulations
  - binary responses
  - Poisson responses
- Revisit the example
- Goals for computational environmetrics

# Goals for Computational Environmetrics

- reproducibility and ease of implementation, particularly for Bayesian methods
- modularity
- comparison and evaluation of models and fitting methods

# GAM and GLMM frameworks

- basic model

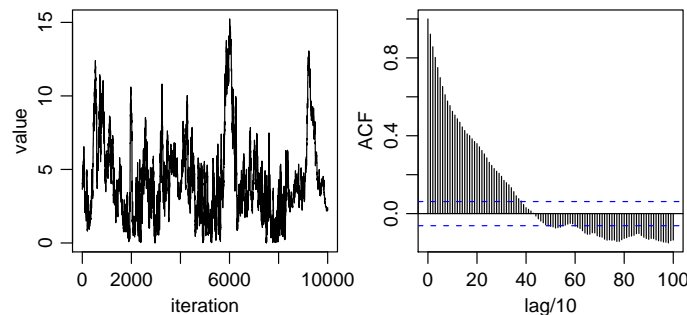
$$Y_i \sim \text{Ber}(p(\mathbf{x}_i, \mathbf{s}_i))$$
$$\text{logit}(p(\mathbf{x}_i, \mathbf{s}_i)) = \mathbf{x}_i^T \boldsymbol{\beta} + g_\theta(\mathbf{s}_i)$$

- basic spatial model for  $\mathbf{g}_\theta^s = (g_\theta(\mathbf{s}_1), \dots, g_\theta(\mathbf{s}_n))$ 
  - GAM:  $g_\theta(\cdot)$  is a two-dimensional smooth term
    - \* basis representation,  $\mathbf{g}_\theta^s = Z\mathbf{u}$
    - \* Gaussian process representation:  $g(\cdot) \sim \text{GP}(\mu(\cdot), C_\theta(\cdot, \cdot)) \Rightarrow \mathbf{g}_\theta^s \sim N(\boldsymbol{\mu}, C_\theta)$
  - GLMM:  $g_\theta(\mathbf{s}_i) = \mathbf{z}_i^T \mathbf{u}$ 
    - \* correlated random effects,  $\mathbf{u} \sim N(0, \Sigma)$



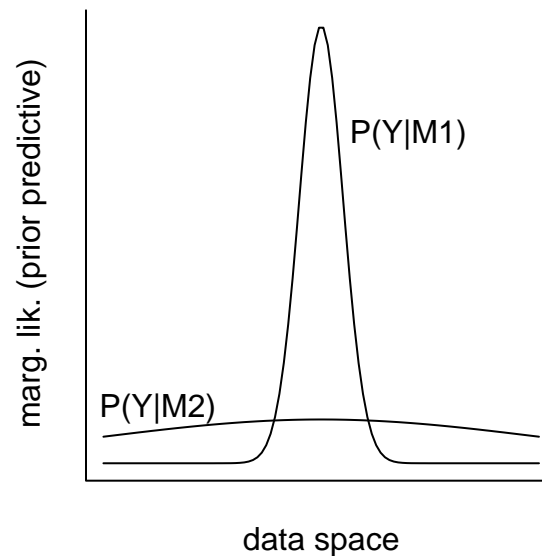
# Difficulties: speed and mixing

- Gaussian responses: closed form marginal likelihood – estimate  $\beta, \theta$
- non-Gaussian: no closed form -> high dimensional estimation – estimate  $\beta, \theta, u$
- Challenges:
  - Classical mixed model: how approximate integral over random effects?
  - Bayesian methods: how perform large matrix calculations and avoid poor mixing?



# Fitting approaches

- penalized likelihood,  $l(\mathbf{y}; \boldsymbol{\beta}, \mathbf{g}_\theta^s) - \lambda J(\mathbf{u})$ 
  - fit by iterative weighted least squares
- Bayesian model for  $(\boldsymbol{\beta}, \theta, \mathbf{u})$ :  $l(\mathbf{y}; \boldsymbol{\beta}, \mathbf{g}_\theta^s) + \log \pi(\boldsymbol{\beta}, \theta, \mathbf{u})$ 
  - fit by MCMC
  - implicit Bayesian penalty on complex spatial functions



# Goals for implementations

- fast computations, avoiding large matrix calculations
- methods that scale reasonably with  $n$
- reasonable fitting of simple risk surfaces we expect to model
- ease of implementation for applied work

# Models and fitting methods considered

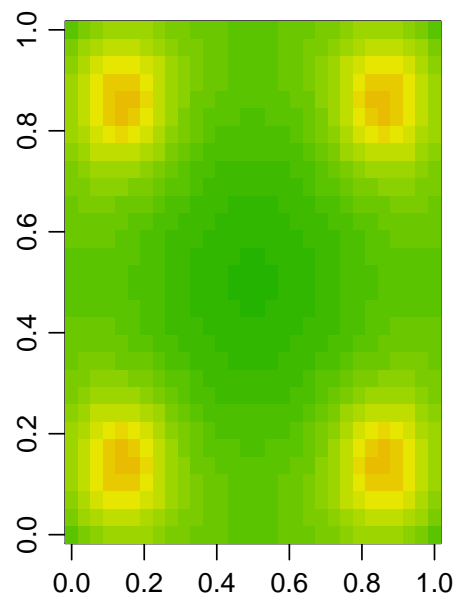
- penalized likelihood based on mixed model with REML smoothing (Kammann and Wand, 2003; Ngo and Wand, 2004) [PL-PQL]
- penalized likelihood with GCV smoothing (Wood, 2001, 2003, 2004) [PL-GCV]
- Bayesian geosadditive model-style radial basis functions fit by MCMC (Zhao and Wand 2004) [B-Geo]
- Bayesian spectral basis representation fit by MCMC using the FFT (Wikle 2002; Paciorek and Ryan, in prep.) [B-SB]
- Bayesian neural network model fit by MCMC (R. Neal) [B-NN]

# Penalized likelihood using GLMM framework with REML [PL-PQL]

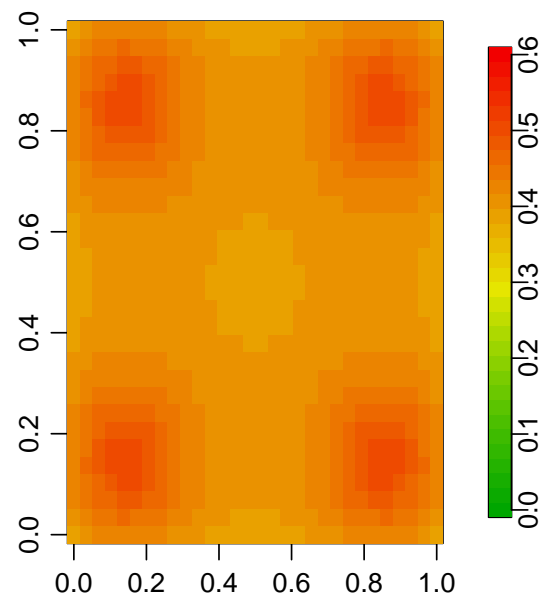
- $\mathbf{g}^s = Z\mathbf{u}$ ,  $Z = \Psi_{nk}\Omega_{kk}^{-\frac{1}{2}}$ ,  $\mathbf{u} \sim N(0, \sigma_u^2)$  - variance component provides complexity penalty
- $\Omega$  contains pairwise spatial covariances between  $k$  knot locations and  $\Psi$  between  $n$  data locations and  $k$  knot locations
- potential covariance functions:
  - thin plate spline generalized covariance function,  $C(\tau) = \tau^2 \log \tau$
  - Matérn correlation function,  $R(\tau) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{2\sqrt{\nu}\tau}{\rho}\right) K_\nu \left(\frac{2\sqrt{\nu}\tau}{\rho}\right)$ , with  $\rho$  and  $\nu$  fixed
- computationally efficient approximation of a Gaussian process representation for  $\mathbf{g}^s$
- PQL approach - IWLS fitting of  $(\beta, \mathbf{u})$  with REML estimation of  $\sigma_u^2$  within the iterations using MM software

# GLMM basis functions

- radial basis functions centered at the knots
- 4 of 64 functions displayed:



TPS



Matérn

# Penalized likelihood using GCV [PL-GCV]

- thin plate spline basis for  $g(\cdot)$
- truncated eigendecomposition of basis matrix increases computational efficiency
- IWLS fitting of  $(\beta, \mathbf{u})$  with GCV estimation of penalty
- easy implementation using the R mgcv library – `gam()`

# Bayesian geoadditive model [B-Geo]

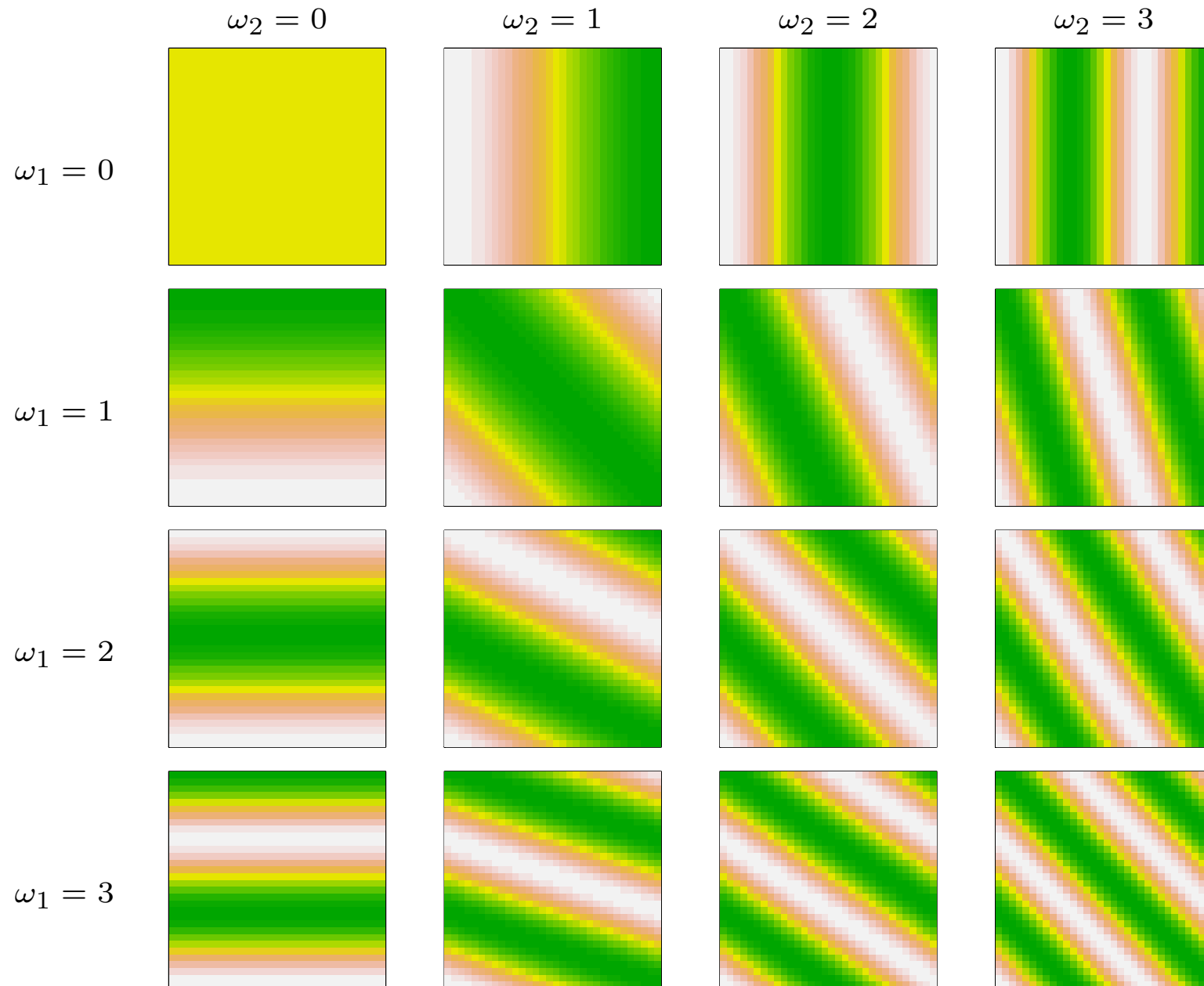
- Bayesian version of GLMM framework already described
  - $\mathbf{g}^s = Z\mathbf{u}$ ,  $Z = \Psi_{nk}\Omega_{kk}^{-\frac{1}{2}}$ ,  $\mathbf{u} \sim N(0, \sigma_u^2)$
  - natural Bayesian complexity penalty through prior on  $\mathbf{u}$
- thin plate spline covariance or Matérn correlation basis construction of  $\Psi$  and  $\Omega$
- MCMC implementation - ensuring mixing is not simple
  - Metropolis-Hastings for  $\mathbf{u}$  using conditional posterior mean and variance based on linearized observations
  - joint proposals for  $\sigma_u^2$  and  $\mathbf{u}$  to ensure that  $\mathbf{u}$  remains compatible with its variance component



# Bayesian spectral basis function model [B-SB]

- computationally efficient basis function construction
- $\mathbf{g}^\# = Z\mathbf{u}$ ,  $\mathbf{g}^s = \sigma P\mathbf{g}^\#$  - piecewise constant gridded surface on  $k$  by  $k$  grid
- $Z$  is the Fourier (spectral) basis and  $Z\mathbf{u}$  is the inverse FFT
- $Z\mathbf{u}$  is approximately a Gaussian process (GP) when...
  - spectral density,  $\pi_\theta(\cdot)$ , of GP covariance function defines  $V(\mathbf{u})$
  - $\mathbf{u} \sim N(0, \text{diag}(\pi_\theta(\boldsymbol{\omega})))$  for Fourier frequencies,  $\boldsymbol{\omega}$

# Bayesian spectral basis functions

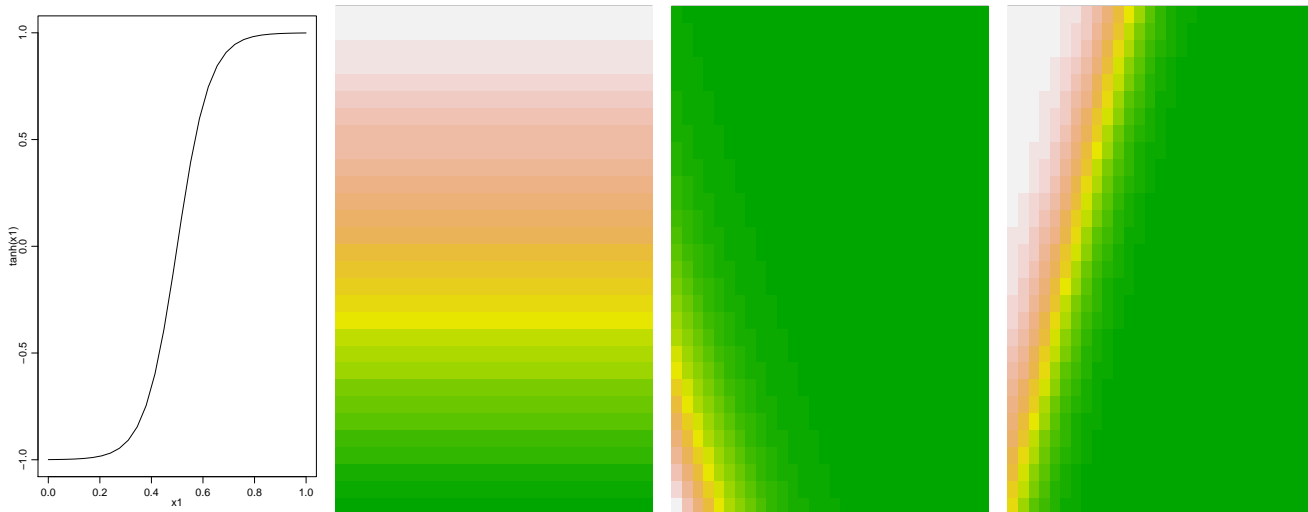


# Comparison with usual GP specification

- usual GP model:  $\mathbf{g}^s \sim N(\boldsymbol{\mu}, C_\theta)$ 
  - $O(n^3)$  fitting:  $|C_\theta|$  and  $C_\theta^{-1}$
- spectral basis uses FFT
  - $O(k^2) \log(k^2)$
  - fast computation and prediction of surface given coefficients
  - a priori independent coefficients give fast computation of prior and help with mixing
  - additional observations are essentially free for a fixed grid

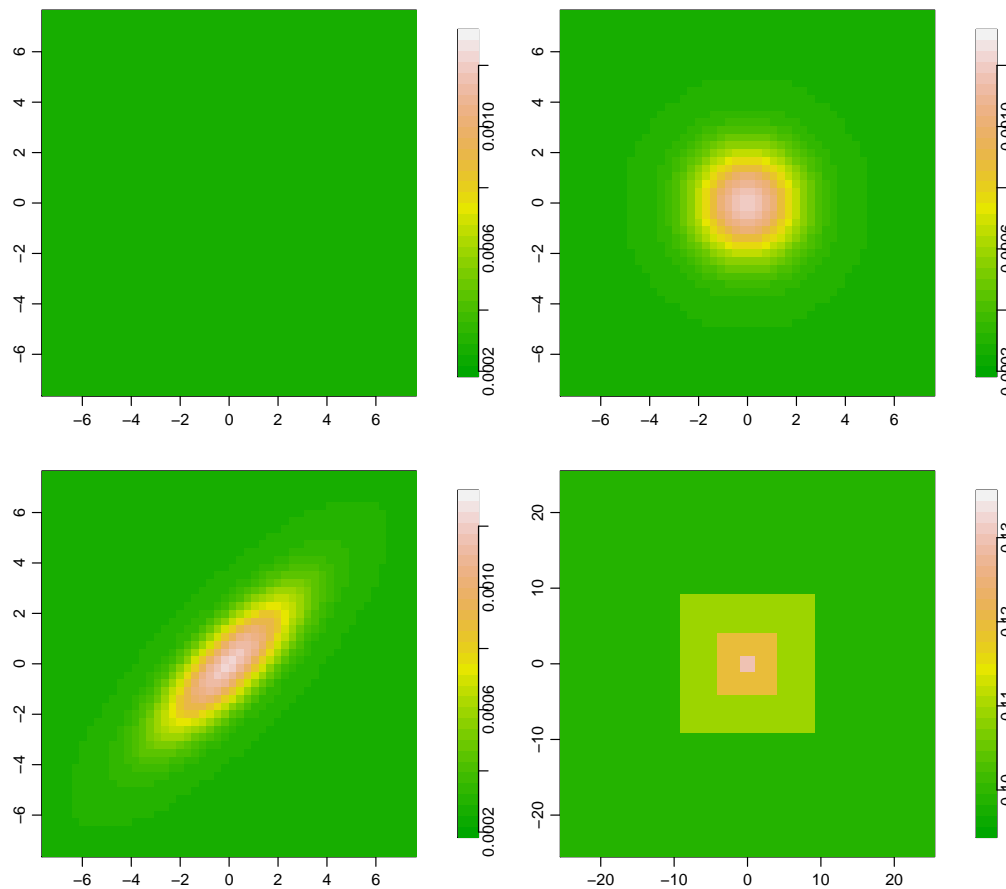
# Bayesian neural network [B-NN]

- multilayer perceptron with one hidden layer gives basis representation:
  - $g(\mathbf{s}_i) = \sum_k \tanh(\phi_k^T \mathbf{s}_i) u_k$
- position and orientation of basis functions change with  $\phi_k$
- implemented with software of R. Neal; somewhat complicated proposal scheme

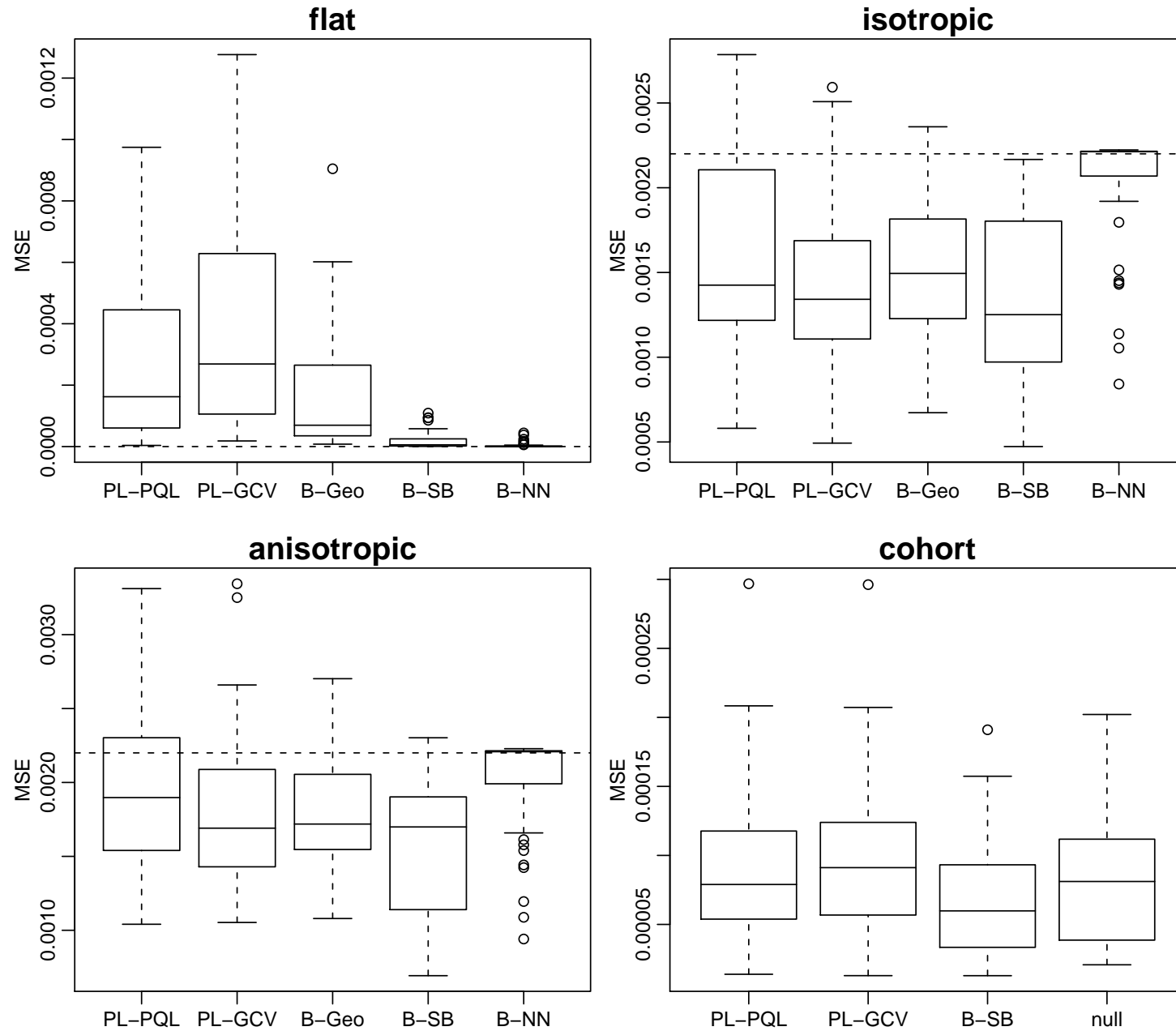


# Simulated datasets

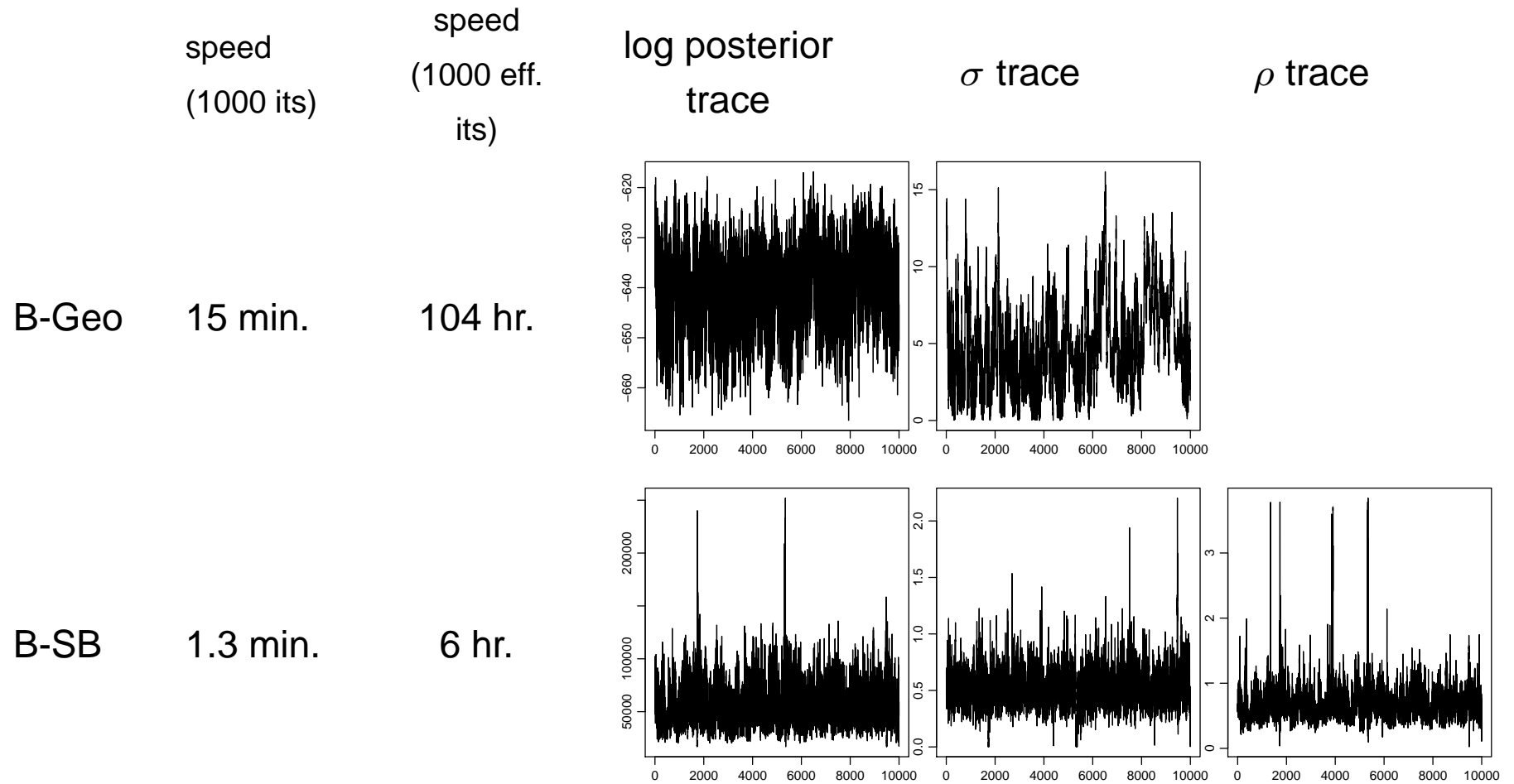
- 3 case-control scenarios:  $n_0 = 1,000$ ;  $n_1 = 200$ ;  $n_{\text{test}} = 2500$  on 50 by 50 grid
- 1 cohort scenario:  $n = 10,000$ ;  $n_{\text{test}} = 2500$  on 50 by 50 grid



# Assessment on 50 simulated datasets



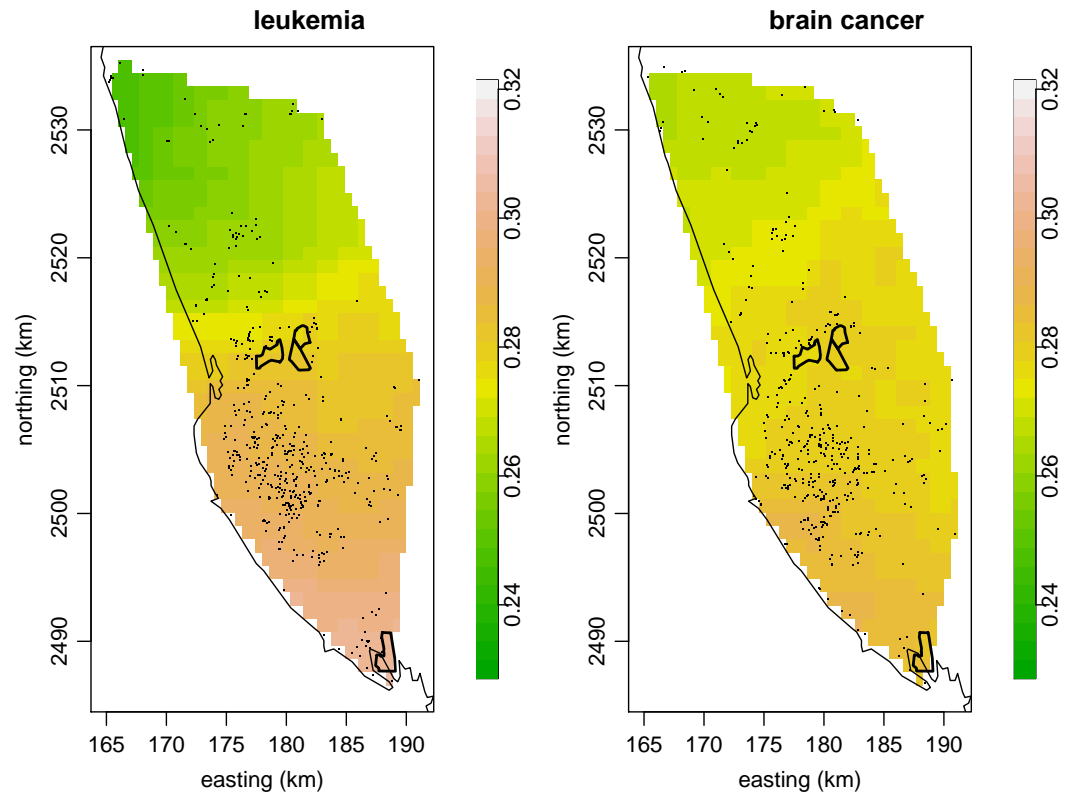
# Mixing and speed of Bayesian methods



# Example revisited - assessment

Summed test deviance  
over 10-fold C-V sets

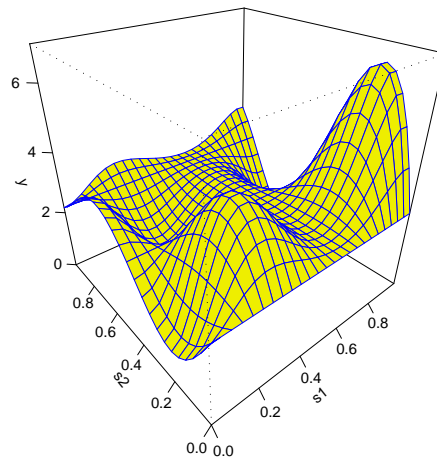
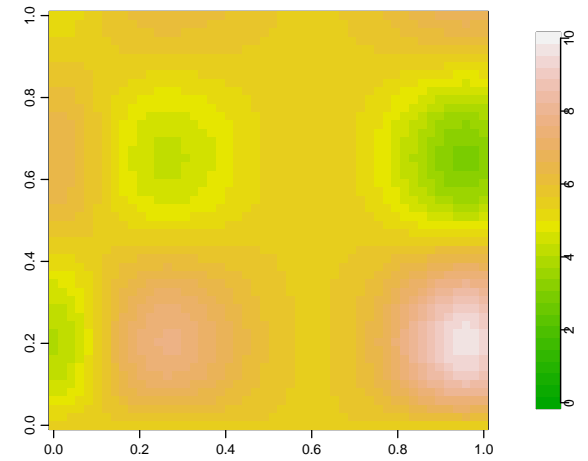
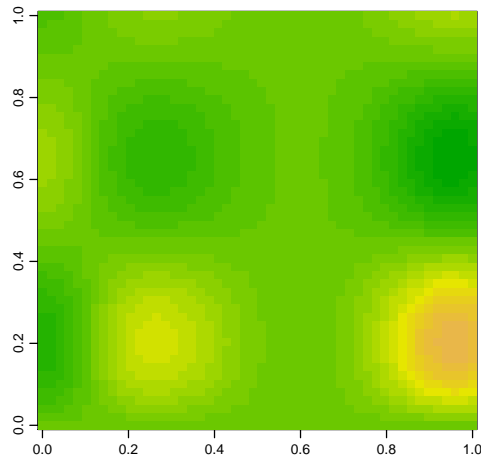
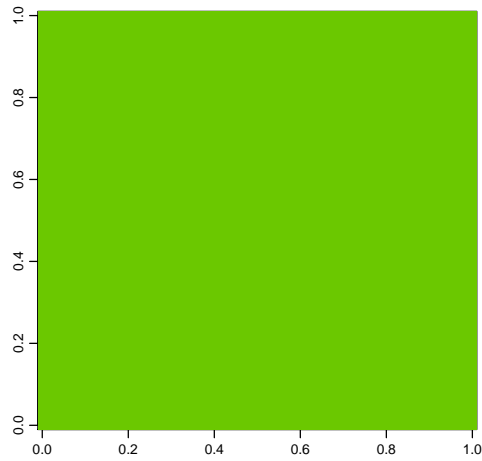
	leukemia	brain cancer
PL-GCV	590.1	529.8
PL-PQL	585.6	529.5
B-Geo	583.3	525.7
B-SB	582.1	525.1
null	581.6	525.5



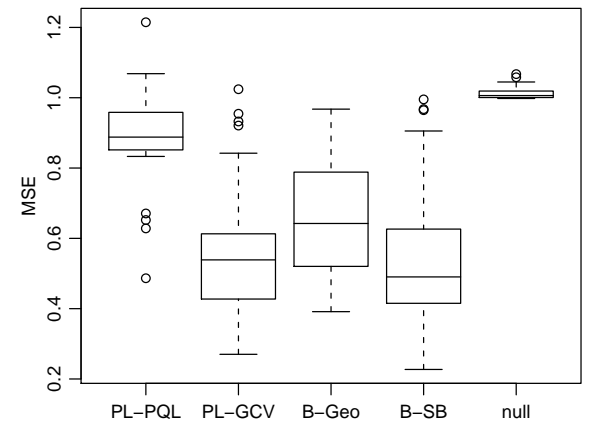
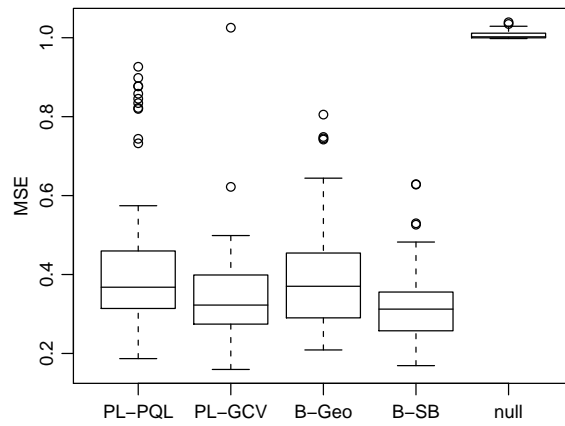
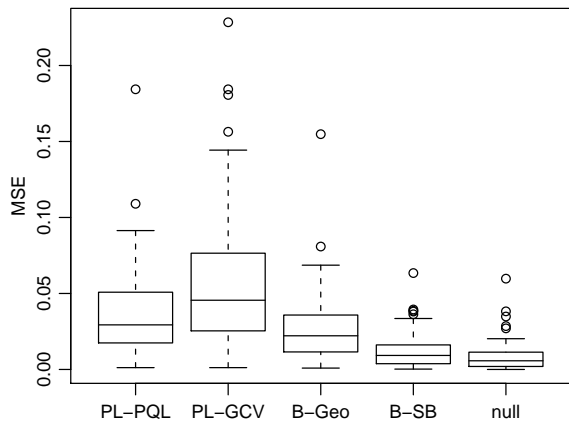
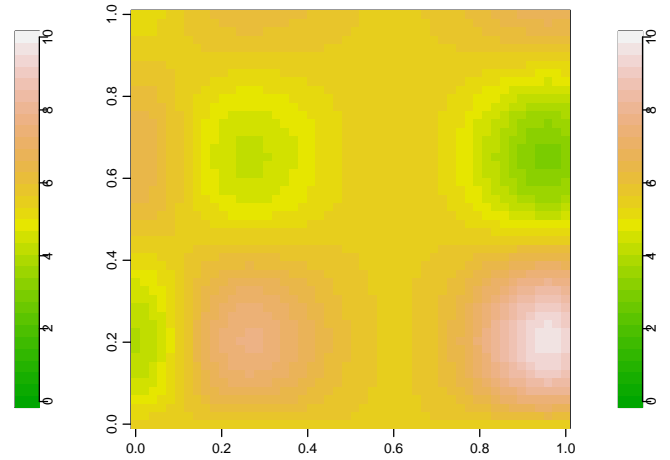
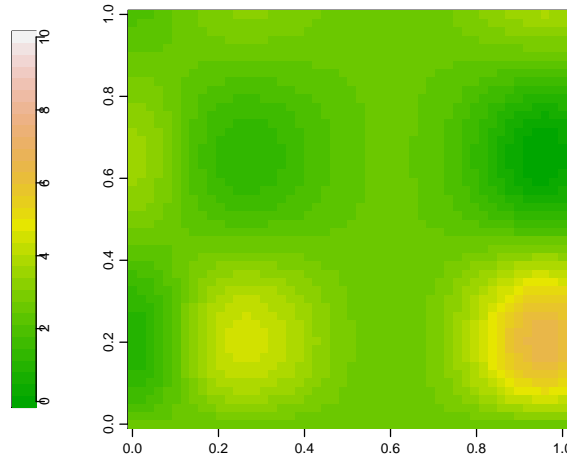
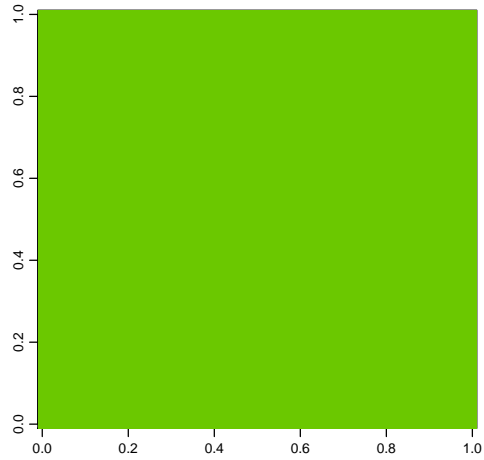


# Simulated count data

- $n = 225$ ,  $n_{\text{test}} = 2500$  on 50 by 50 grid



# Assessment on count simulations



# Methodology lessons

- Effective process parameterization allows for faster Bayesian estimation
  - effective for spatial models with thousands of observations
- Natural Bayesian complexity penalty works well; other automatic criteria appear to overfit
- R code for spectral basis model to ease implementation
- Power is an issue with binary observations
- Results hold for count data
- Spectral basis could provide a modular building block in hierarchical models

# Suggestions for computational environmetrics

- reproducibility
  - requires code and detailed description (supplemental material/web)
  - standard computing environment (R) helps
  - enabling reproducible MCMC (beyond BUGS)- class structures, templates, and proposal functions for R
- modularity
  - spectral basis as modular component
- comparison of methods
  - rare
  - difficult without reproducibility, particularly with Bayesian methods

# Methodological future work

- Importance of basis functions vs. speed/mixing in MCMC vs. penalty estimation method in determining fitting success
  - Why don't automatic criteria for penalized likelihood work as well?
  - Importance of fitting variance and spatial range parameters
    - small-sample results (consider effective basis functions) vs. asymptotics (Zhang, 2004)
- Simple approaches for testing necessity of spatial term
- Other process parameterizations allowing fast Bayesian estimation:
  - Simple prior structures for wavelet basis coefficient (co)variances?