

Statistical Modelling Approaches to Disease Mapping

Peter J Diggle

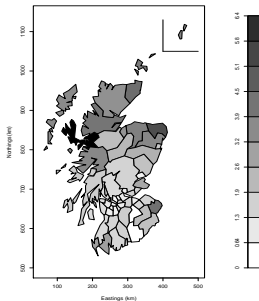
Lancaster University and University of Liverpool



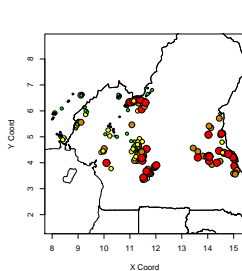
UNIVERSITY OF
LIVERPOOL

INSTITUTE OF INFECTION
AND GLOBAL HEALTH

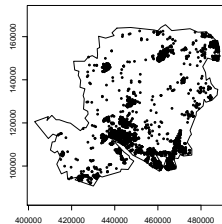
Spatial statistics according to Cressie (1991)



Lattice data

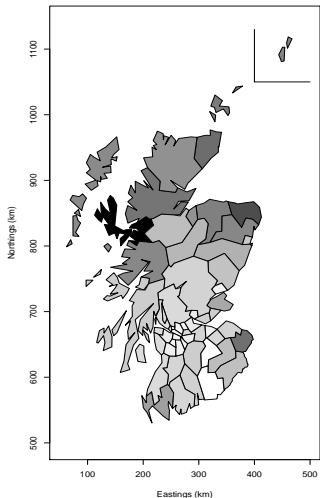


Geostatistics



Point patterns

Lattice data: Scottish lip cancer incidence



Data: county-level incidences

$Y_i : i = 1, \dots, n$

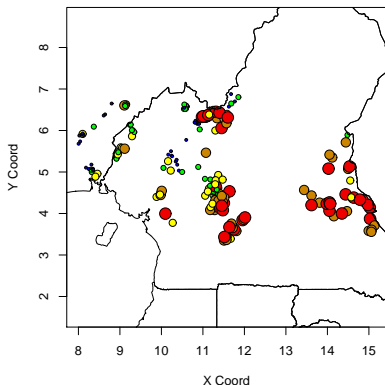
Model: Markov random field:

$[Y_i | \{Y_j : j \neq i\}] : i = 1, \dots, n$

- risks in near-neighbouring counties are positively correlated
- incidences Y_i are noisy versions of $\text{risk} \times \text{population}$

Scientific interest confined to specified set of counties?

Geostatistics: Loa loa prevalence in Cameroon



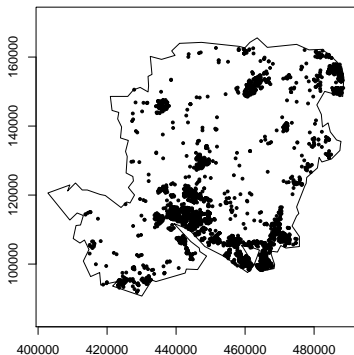
Data: empirical prevalences Y_i at sample locations $x_i : i = 1, \dots, n$

Model: spatially continuous stochastic process, $S(x) : x \in \mathbb{R}^2$

- correlation between $S(u)$ and $S(v)$ specified as a function of distance between u and v
- $Y_i | S(x_i) \sim \text{Binomial}$

Scientific interest extends to $S(x)$ at non-sampled locations

Point pattern: gastro-enteric illness in Hampshire



Data: outcomes (x_i, t_i) are locations and dates of calls to NHS Direct recorded as “vomiting and/or diarrhoea”

Model: $(x_i, t_i) : i = 1, 2, \dots$ a stochastic point process

- intensity $\lambda(x, t)$
- successive cases independent?

Scientific interest is in locations themselves

Context

- region of interest A
- disease risk $\rho(x) : x \in A$
- data relating to variation in disease prevalence over A

Objective

- estimate $\rho(x)$?
- calculate $P\{\rho(x) > c | \text{data}\}$?

The answer to any prediction problem is a probability distribution

Peter McCullagh, FRS

Markov Random Field (MRF) models

(Besag, 1974; Rue and Held, 2005)

- **Random variables** $S = (S_1, \dots, S_n)$
- **Joint distribution** $[S]$ fully specified by **full conditionals**,

$$[S_i | \{S_j : j \neq i\}] : i = 1, \dots, n$$

- **Neighbourhood** of i is $\mathcal{N}(i) \subset \{1, 2, \dots, n\}$

$$[S_i | \{S_j : j \neq i\}] = [S_i | S_j : j \in \mathcal{N}(i)] : i = 1, \dots, n$$

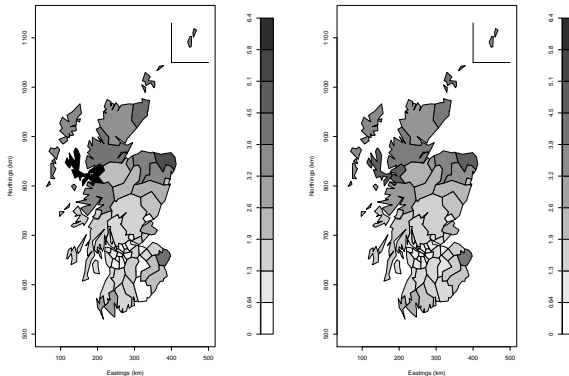
- latent Gaussian MRF $\mathbf{S} = (\mathbf{S}_1, \dots, \mathbf{S}_n)$,

$$\mathbf{S}_i | \{\mathbf{S}_j : j \neq i\} \sim \mathcal{N}(\bar{\mathbf{S}}_i, \tau^2 / m_i)$$

- conditionally independent $\mathbf{Y}_i | \mathbf{S} \sim \text{Poiss}(\mathbf{z}_i' \beta + \gamma \mathbf{S}_i)$
- risk map: $\mathbb{E}[\mathbf{S}_i | \mathbf{Y}]$

Besag, York and Mollié, 1991

Raw and spatially smoothed relative risk estimates for lip cancer in 56 Scottish counties



Wakefield (2007)

Limitations of MRF models for spatial data

MRF's are just multivariate probability distributions

- parameterised in a way that has a spatial interpretation
- but specific to a fixed set of locations x_1, \dots, x_n

Neighbourhood specification can be problematic

- natural hierarchy of models on regular lattices
- not so for irregular lattices
- and arguably un-natural for spatially aggregated data,

$$Y_i = \int_{A_i} Y(x) dx$$

Geostatistical models

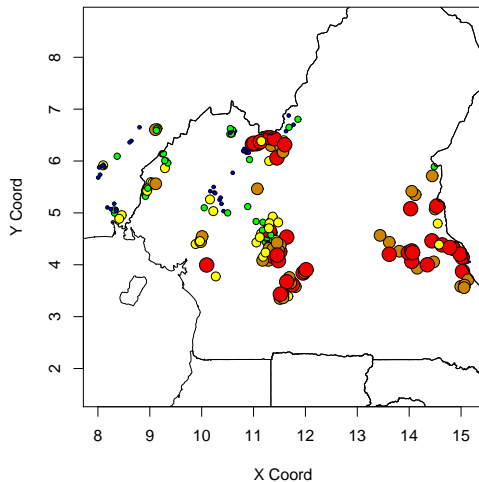
(Diggle and Ribeiro, 2007; Chilès and Delfiner, 2012)

- **Stochastic process** $S(\mathbf{x}) : \mathbf{x} \in \mathbf{A} \subset \mathbb{R}^2$
- **Data** $\{(Y_i, \mathbf{x}_i) : i = 1, \dots, n\}$
- **Stationary Gaussian model**

$$E[S(\mathbf{x}) = 0] \quad \text{Cov}\{S(\mathbf{x}), S(\mathbf{x} - \mathbf{u})\} = \sigma^2 \rho(\mathbf{u})$$

$$[Y|S] = [Y_1|S(\mathbf{x}_1)] \dots [Y_n|S(\mathbf{x}_n)]$$

A geostatistical data-set: Loa loa prevalence surveys



- **Latent spatially correlated process**

$$\mathbf{S}(\mathbf{x}) \sim \text{SGP}\{0, \sigma^2, \rho(\mathbf{u})\}$$
$$\rho(\mathbf{u}) = \exp(-|\mathbf{u}|/\phi)$$

- **Linear predictor (regression model)**

$$\mathbf{d}(\mathbf{x}) = \text{environmental variables at location } \mathbf{x}$$
$$\eta(\mathbf{x}) = \mathbf{d}(\mathbf{x})'\beta + \mathbf{S}(\mathbf{x})$$
$$\mathbf{p}(\mathbf{x}) = \log[\eta(\mathbf{x})/\{1 - \eta(\mathbf{x})\}]$$

- **Conditional distribution for positive proportion \mathbf{Y}_i/n_i**

$$\mathbf{Y}_i | \mathbf{S}(\cdot) \sim \text{Bin}\{n_i, \mathbf{p}(\mathbf{x}_i)\} \text{ (binomial sampling)}$$

Probabilistic exceedance map for Cameroon (Diggle et al, 2007)

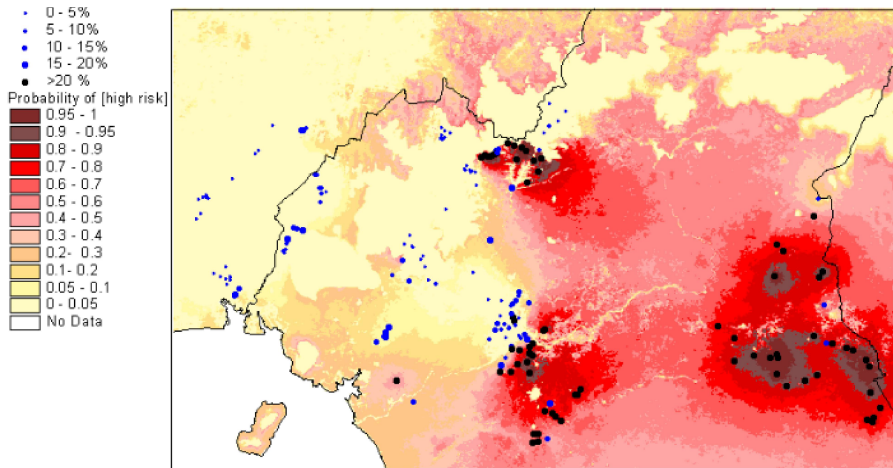


Figure 6: 'PCM for [high risk] in Cameroon based on 'ERM with ground truth data.

Point process models: log-Gaussian Cox process (Møller, Syversveen and Waagepetersen, 1998)

- **Stochastic process** $S(\mathbf{x}) : \mathbf{x} \in \mathbf{A} \subset \mathbb{R}^2$
- **Data** $\mathcal{X} = \{\mathbf{x}_i : i = 1, \dots, n\}$
- **Stationary Gaussian model**

$$E[S(\mathbf{x}) = 0] \quad \text{Cov}\{S(\mathbf{x}), S(\mathbf{x} - \mathbf{u})\} = \sigma^2 \rho(\mathbf{u})$$

$$[\mathcal{X} | S] = \text{Poisson process, intensity } \Lambda(\mathbf{x}) = \exp\{S(\mathbf{x})\}$$

Real-time surveillance: spatio-temporal point process (Diggle, Rowlingson and Su, 2005)

Ascertainment and
Enhancement of
Gastroenteric
Infection
Surveillance
Statistics

- largely sporadic incidence pattern
- concentration in population centres
- occasional “clusters” of cases

Can spatial statistical modelling enable earlier detection of “clusters”?

$$\text{intensity} = \text{expected} \times \text{unexpected}$$

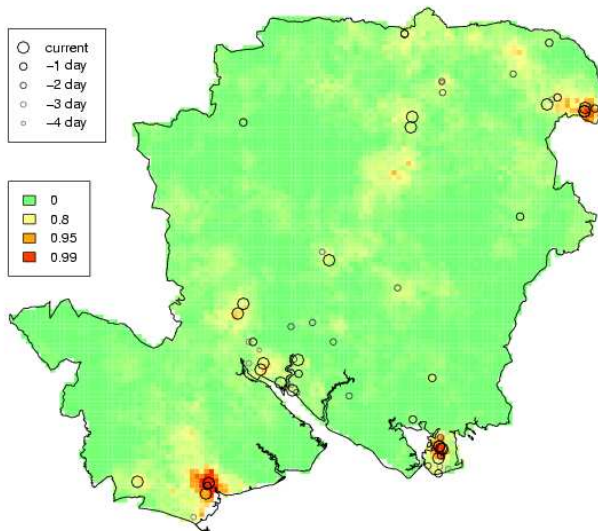
$$\Lambda(x, t) = \lambda_0(x) \times \mu_0(t) \times R(x, t)$$

Objective: use incident data up to time t to construct predictive distribution for current “anomaly” surface, $R(x, t)$

Model

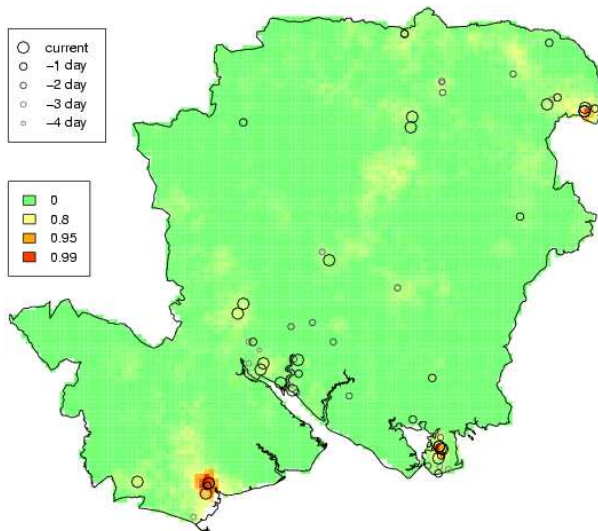
- spatio-temporal point process \mathcal{P}
- $\log R(x, t) \sim$ latent Gaussian process
- $\mathcal{P}|R \sim$ Poisson process

Spatial prediction: 6 March 2003



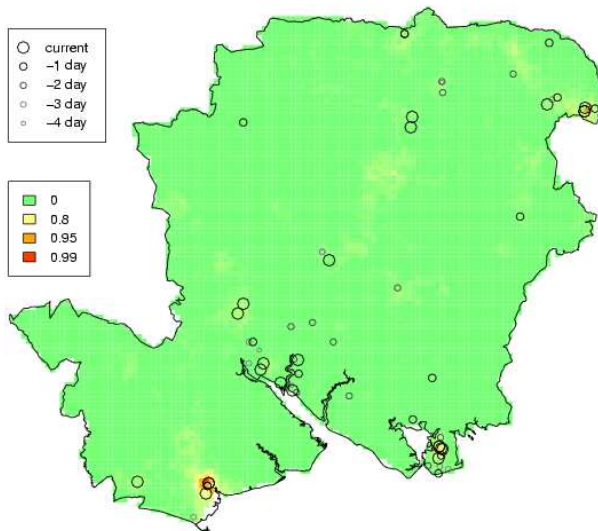
$c = 2$

Spatial prediction: 6 March 2003



$c = 4$

Spatial prediction: 6 March 2003



$c = 8$

Synthesis

Diggle, Moraga, Rowlingson, and Taylor, 2013)

S = state of nature
Y = **all** relevant data
T = $\mathcal{F}(S)$ = target for prediction

Model: $[S, Y] = [S][Y|S]$
Prediction: $[S, Y] \Rightarrow [S|Y] \Rightarrow [T|Y]$

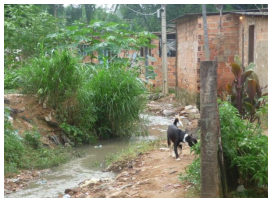
Pau da Lima, Salvador, Brazil



Pau da Lima, Salvador, Brazil



Leptospirosis cohort study: Pau da Lima



- **subjects** i at locations x_i , **blood-samples** taken at times $t_{ij} \approx 0, 6, 12, 18, 24$ months
- **sero-conversion** defined as change from zero to positive, or at least four-fold increase in concentration
- **data** consist of:
 - $Y_{ij} = 0/1 : j = 1, 2, 3, 4$ (seroconversion no/yes)
 - $r_i(t)$ known and hypothesised risk-factors

Leptospirosis cohort study: analysing the data

Longitudinal data, binary outcome \Rightarrow standard problem?

id	Follow-up				Age
	1	2	3	4	
1	0	0	1	0	57
2	0	0	0	0	34
3	0	0	1	X	38
4	1	1	1	0	28
.
.
.
950	0	1	0	1	40

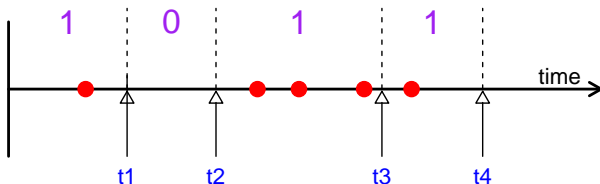
Logistic regression for binary response,

$$\log\{p_{it}/(1 - p_{it})\} = \alpha + \beta \times \text{age}$$

Need to account for correlation amongst repeated outcomes on same individual

- generalized estimating equations
- generalized linear mixed models
- ...

Leptospirosis cohort study: analysing the problem



- **infection events** on each individual form a point process with time-varying intensity, $\Lambda_i(t)$
- **follow-up times** partially censor the point process record
- **reduction to binary data** represents additional censoring

Leptospirosis cohort study: model formulation

Data: $Y_{it} = 0/1 \quad t = 1, 2, 3, 4 \quad i = 1, 2, \dots, n$

- $Y_{it} = 1 \Leftrightarrow$ at least one infection event
- model infection events as person-specific, inhomogeneous Cox processes,

$$\Lambda_i(t) = \exp\{r_i(t)' \beta + U_i + S(x_i)\}$$

$$P(Y_{it} = 1 | \Lambda_i(\cdot)) = 1 - \exp \left\{ - \int_{t_{i,j-1}}^{t_{ij}} \Lambda_i(u) du \right\}$$

- **The likelihood principle**

Two data-sets x and y that generate identical likelihood functions are equivalent as evidence

- **The law of likelihood**

If $H_A \Rightarrow p_A(x)$ and $H_B \Rightarrow p_B(x)$, then data x constitutes evidence in favour of A over B iff $p_A(x) > p_B(x)$, and the likelihood ratio, $p_A(x)/p_B(x)$ measures the strength of the evidence

Inference: what's the question? (Royall, 1997)

- **Bayesian**

What should I believe?

- **Decision-theoretic**

What should I do?

- **Classical:**

What do the data tell me?

Acknowledgements

CHICAS, Lancaster University : Paula Moraga, Barry Rowlingson, Ben Taylor

APOC Madeleine Thomson, Hans Remme, Honorat Zoure, ...

Yale University/Fiocruz, Brazil: Federico Costa, Jose Hagan, Albert Ko

MRC: Methodology Research Grant G0902153

References

- Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems (with Discussion). *Journal of the Royal Statistical Society B* **36**, 192–225.
- Besag, J., York, J. and Mollié, A. (1991). Bayesian image restoration, with two applications in spatial statistics (with Discussion). *Annals of the Institute of Statistical Mathematics*, **43**, 1–59.
- Chilès, J-P and Delfiner, P. (2012). *Geostatistics* (second edition). Hoboken: Wiley.
- Cressie, N.A.C. (1991). *Statistics for Spatial Data*. New York: Wiley.
- Diggle, P.J., Moraga, P., Rowlingson, B. and Taylor, B. (2013). Spatial and spatio-temporal log-Gaussian Cox processes: extending the geostatistical paradigm. *Statistical Science*, **28**, 542–563.
- Diggle, P.J. and Ribeiro, P.J. (2007). *Model-based Geostatistics*. New York: Springer.
- Diggle, P., Rowlingson, B. and Su, T. (2005). Point process methodology for on-line spatio-temporal disease surveillance. *Environmetrics*, **16**, 423–34.
- Diggle, P.J., Thomson, M.C., Christensen, O.F., Rowlingson, B., Obsomer, V., Gardon, J., Wanji, S., Takougang, I., Enyong, P., Kamgno, J., Remme, H., Boussinesq, M. and Molyneux, D.H. (2007). Spatial modelling and prediction of Loa loa risk: decision making under uncertainty. *Annals of Tropical Medicine and Parasitology*, **101**, 499–509.
- Møller, J., Syversveen, A. and Waagepetersen, R. (1998). Log Gaussian Cox processes. *Scandinavian Journal of Statistics*, **25**, 451–82.
- Royall, R. (1997). *Statistical Evidence: a likelihood paradigm*. London: Chapman and Hall.
- Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*. London: CRC Press.
- Wakefield, J. (2007). Disease mapping and spatial regression with count data. *Biostatistics*, **8**, 158–183.