Statistical Modelling Approaches to Disease Mapping

Peter J Diggle

Lancaster University and University of Liverpool





Spatial statistics according to Cressie (1991)



Lattice data

Geostatistics

Point patterns

Cressie, N.A.C. (1991). Statistics for Spatial Data. Wiley.

Lattice data: Scottish lip cancer incidence



- risks in near-neighbouring counties are positively correlated
- incidences \textbf{Y}_i are noisy versions of risk \times population

Scientific interest confined to specified set of counties?



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

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



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, ... a stochastic point process

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

Scientific interest is in locations themselves

Context

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

Objective

- estimate ρ(x) ?
- calculate P{ρ(x) > c|data}?

The answer to any prediction problem is a probability distribution Peter McCullagh, FRS

Markov Random Field (MRF) models

- Random variables S = (S₁, ..., S_n)
- Joint distribution [S] fully specified by full conditionals,

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

• Neighbourhood of i is $\mathcal{N}(\mathsf{i}) \subset \{1,2,...,\mathsf{n}\}$

 $[\textbf{S}_i|\{\textbf{S}_j:j\neq i\}]=[\textbf{S}_i|\textbf{S}_j:j\in\mathcal{N}(i)]:i=1,...,n$

Hierarchical Poisson/Gaussian MRF

• latent Gaussian MRF $S = (S_1, ..., S_n)$,

 $S_i|\{S_j:j\neq i\}\sim \mathrm{N}(\bar{S}_i,\tau^2/m_i)$

• conditionally independent $Y_i | S \sim Poiss(z'_i \beta + \gamma S_i)$

• risk map: **E[S**_i|**Y**]

Besag, York and Mollié, 1991

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





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 x₁, ..., 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,

$$\mathbf{Y}_{i} = \int_{\mathbf{A}_{i}} \mathbf{Y}(\mathbf{x}) d\mathbf{x}$$

Geostatistical models

• Stochastic process $S(x): x \in A \subset \mathbb{R}^2$

• Data
$$\{(Y_i, x_i) : i = 1, ..., n\}$$

• Stationary Gaussian model

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

 $[\boldsymbol{Y}|\boldsymbol{S}] = [\boldsymbol{Y}_1|\boldsymbol{S}(\boldsymbol{x}_1)]...[\boldsymbol{Y}_n|\boldsymbol{S}(\boldsymbol{x}_n)]$

A geostatistical data-set: Loa loa prevalence surveys



X Coord

コントロン キョント ヨー ろくの

Loa loa: generalised linear model

Latent spatially correlated process

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

• Linear predictor (regression model)

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

• Conditional distribution for positive proportion Y_i/n_i $Y_i|S(\cdot) \sim Bin\{n_i, p(x_i)\}$ (binomial sampling)

▲□▶ ▲□▶ ▲目▶ ▲目▶ - 目 - のへの

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



Figure 6: PCM for /high risk/ in Cameroon based on ERMr with ground truth data.

Point process models (log-Gaussian Cox processes)

• Stochastic process $S(x) : x \in A \subset \mathbb{R}^2$

• Data
$$\mathcal{X} = \{x_i : i = 1, ..., n\}$$

• Stationary Gaussian model

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

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

Real-time spatial surveillance: spatio-temporal point process

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"?

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

Model

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

Spatial prediction: 6 March 2003



Spatial prediction: 6 March 2003



◆□> ◆□> ◆目> ◆目> ◆目> 目 のへで

Spatial prediction: 6 March 2003



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

$\begin{array}{ll} \mbox{Model:} & [S,Y] = [S][Y|S] \\ \mbox{Prediction:} & [S,Y] \Rightarrow [S|Y] \Rightarrow [T|Y] \end{array}$

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* (to appear)



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)
 - $\bullet\ r_i(t)$ known and hypothesised risk-factors

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	Х	38
4	1	1	1	0	28
•	•	•	•	•	•
•	•	•	•	•	•
•	•	•	•	•	•
950	0	1	0	1	40

Logistic regression for binary response,

• ...

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

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

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

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

$$\Lambda_{i}(t) = \exp\{r_{i}(t)'\beta + U_{i} + S(x_{i})\}$$

$$\mathbf{P}(\mathbf{Y}_{it} = 1 | \mathbf{\Lambda}_i(\cdot)\} = 1 - \exp\left\{-\int_{t_{i,j-1}}^{t_{ij}} \mathbf{\Lambda}_i(\mathbf{u}) d\mathbf{u}\right\}$$

• 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

Bayesian

What should I believe?

• Decision-theoretic

What should I do?

• Classical:

What do the data tell me?

Royall, R. (1997). Statistical Evidence: a likelihood paradigm. London: Chapman and Hall.

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ の 0 0

MRC: Methodology Research Grant G0902153