Statistical modelling and inference for spatio-temporal disease processes

Martin Hazelton¹

Statistics and Bioinformatics Group Institute of Fundamental Sciences, Massey University

24 October 2012



Presenter: m.hazelton@massey.ac.nz

Disease, Prediction and Chance



- Spread of infectious disease is a complex process typically involving a plethora of factors.
- Some are explicable/predictable.
- Others are essentially random.
- Plenty of work for statisticians deciding which is which.



Spatial Patterns ... of Disease?



- Complete spatial randomness
- Locations of gastroenteritis cases
- Mystery spatial point pattern



Spatial Patterns ... of Disease?



- Complete spatial randomness (B)
- Locations of gastroenteritis cases (A)
- Mystery spatial point pattern (C)



Measuring Spatial Risk

- To understand spatial patterns of disease we must adjust for underlying population density.
- To this end, Bithell (1990) defined the relative risk function:

$$r(\boldsymbol{z}) = rac{f(\boldsymbol{z})}{g(\boldsymbol{z})}$$
 $\boldsymbol{z} \in \mathcal{R}$

where f is density of cases and g density of controls.

• Usual to work on log-scale:

$$\rho(\boldsymbol{z}) = \log\{r(\boldsymbol{z})\} = \log\{f(\boldsymbol{z})\} - \log\{g(\boldsymbol{z})\}$$

Bithell J.F. (1990). An application of density estimation to geographical epidemiology. *Statistics in Medicine* **9**:691–701.

Kernel Estimation

• Straightforward approach to estimation of *r* is to replacing unknown densities by kernel estimates thereof:

$$\hat{r}(\boldsymbol{z}) = rac{\hat{f}(\boldsymbol{z})}{\hat{g}(\boldsymbol{z})} \qquad \boldsymbol{z} \in \mathcal{R}$$

• Kernel density estimate constructed from bivariate data x_1, \ldots, x_n :

$$\hat{f}(\boldsymbol{z}) = n^{-1} \sum_{i=1}^{n} K_h(\boldsymbol{z} - \boldsymbol{x}_i)$$

- $K_h(\mathbf{z}) = h^{-2}K(\mathbf{z}/h)$, with kernel K being spherically symmetric;
- h is the bandwidth, controlling degree of smoothing.



An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





An example using the gastroenteritis data





Small bandwidth (undersmoothing)





Small bandwidth (undersmoothing)





Small bandwidth (undersmoothing)





Small bandwidth (undersmoothing)





Large bandwidth (oversmoothing)





Large bandwidth (oversmoothing)





Large bandwidth (oversmoothing)





Large bandwidth (oversmoothing)





Example: Cancer of the Larynx in South Lancashire

- Data are geographical coordinates for 58 cases of cancer of the larynx, and 978 controls (cases of lung cancer).
- Data collected between 1974 and 1983 by Chorley and South Ribble Health Authority in Lancashire, England.



The Effect of Bandwidth Selection

For the log-relative risk function





Research on Bandwidth Selection

 Bandwidth selection is a challenging problem for both theoretical and practical reasons.

- Typical inhomogeneity of populations.
- Want stable estimates of risk even where data are sparse.
- Asymptotic theory very delicate.
- Recent progress on spatially adaptive smoothing regimens Davies & Hazelton (2010).

Davies, T.M. and Hazelton, M.L. (2010). Adaptive kernel estimation of spatial relative risk. *Statistics in Medicine* **29**, 2423–2437.



Estimation of case density





Estimation of case density





Estimation of case density





Estimation of case density





Estimation of case density





Estimation of case density





Estimation of case density





Estimation of log-relative risk surface





Boundary Bias

- Disease data collected over finite region.
- For points near edge, some of kernel weight spills over boundary and is lost.
- Has significant impact on kernel estimates.
- Theoretical analysis is intricate.





Foundations of Asymptotic Analysis of Boundary Bias





IDReC Symposium

15/31

Asymptotics with Adaptive Boundary Kernel Or ... Why Postdocs are Great

$$\begin{split} \frac{\operatorname{bias}(\widehat{f}(\mathbf{z}))}{h_0^2} &= \frac{[D^{10}f(\mathbf{z})]^2}{4f(\mathbf{z})} \left[\frac{b_0}{q} \left(a_{40}^{(20)} + 2a_{31}^{(11)} + a_{22}^{(02)} + 7a_{30}^{(10)} + 7a_{21}^{(01)} + 8a_{20} \right) \\ &\quad + \frac{b_1}{q} \left(a_{50}^{(20)} + 2a_{41}^{(11)} + a_{32}^{(02)} + 9a_{40}^{(10)} + 9a_{31}^{(01)} + 15a_{30} \right) \\ &\quad + \frac{b_2}{q} \left(a_{41}^{(20)} + 2a_{32}^{(12)} + a_{23}^{(02)} + 9a_{31}^{(10)} + 9a_{22}^{(01)} + 15a_{21} \right) \right] \\ &\quad + \frac{D^{10}f(\mathbf{z})D^{01}f(\mathbf{z})}{4f(\mathbf{z})} \left[\frac{b_0}{q} \left(2a_{31}^{(20)} + 4a_{22}^{(11)} + 2a_{13}^{(02)} + 14a_{21}^{(10)} + 14a_{12}^{(01)} + 16a_{11} \right) \\ &\quad + \frac{b_1}{q} \left(2a_{41}^{(20)} + 4a_{22}^{(11)} + 2a_{23}^{(02)} + 18a_{31}^{(10)} + 18a_{22}^{(01)} + 30a_{21} \right) \right] \\ &\quad + \frac{b_2}{q} \left(2a_{32}^{(20)} + 4a_{23}^{(11)} + 2a_{14}^{(02)} + 18a_{13}^{(10)} + 18a_{13}^{(01)} + 30a_{12} \right) \right] \\ &\quad + \frac{b_2}{q} \left(a_{22}^{(20)} + 2a_{13}^{(11)} + a_{04}^{(02)} + 7a_{12}^{(10)} + 7a_{03}^{(01)} + 8a_{02} \right) \\ &\quad + \frac{b_1}{q} \left(a_{32}^{(20)} + 2a_{23}^{(11)} + a_{04}^{(20)} + 9a_{13}^{(10)} + 9a_{04}^{(11)} + 15a_{12} \right) \\ &\quad + \frac{b_2}{q} \left(a_{23}^{(20)} + 2a_{23}^{(11)} + a_{05}^{(2)} + 9a_{13}^{(10)} + 9a_{04}^{(11)} + 15a_{12} \right) \\ &\quad + \frac{b_2}{q} \left(a_{23}^{(10)} + 2a_{21}^{(11)} + a_{05}^{(2)} + 9a_{13}^{(10)} + 9a_{04}^{(11)} + 15a_{03} \right) \right] \\ &\quad + \frac{D^{20}f(\mathbf{z})}{4} \left[\frac{b_0}{q} \left(2a_{30}^{(10)} + 2a_{21}^{(11)} + 8a_{20} \right) + \frac{b_1}{q} \left(2a_{40}^{(10)} + 2a_{31}^{(11)} + 10a_{30} \right) \\ &\quad + \frac{b_2}{2} \left(2a_{31}^{(10)} + 2a_{21}^{(11)} + 8a_{20} \right) + \frac{b_1}{q} \left(2a_{40}^{(10)} + 2a_{31}^{(11)} + 10a_{30} \right) \\ &\quad + \frac{b_2}{a} \left(2a_{31}^{(10)} + 2a_{21}^{(11)} + 10a_{21} \right) \right] \end{aligned}$$

IDReC Symposium

16/31

Infectious Disease Resear

Computer Implementation

- Important practical problems in spatial epidemiology can quickly lead to fascinating theoretical work.
- Need to make this stuff available to users.
- Released R package sparr. See Davies, Hazelton & Marshall (2011).

Davies, T.M., Hazelton, M.L. and Marshall, J.C. (2011). sparr: Analyzing spatial relative risk using fixed and adaptive kernel density estimation in R. *Journal of Statistical Software* **39**, 1–14.



Screenshot of R Package Sparr



Statistics and Bioinformatics Group Research

Methodological research with applications to infectious disease

- Spatial statistics
- Markov chain Monte Carlo methods
- Simulation based inference
- Smoothing methods
- Semiparametric modelling
- Statistical software development
- Variety of applications



Some Members of the Research Group

Postgraduate students



Sarojinie Fernando PhD (recently submitted)



Kate Richards PhD (year 2)



Lyndal Henden Honours



PhD (year 3) Khair Jones PhD (year 1)

Brigid Betz-Stablein

Staff with interest in infectious disease applications

Prof Martin Hazelton A/Prof Geoff Jones Dr Chris Jewell Dr Jonathan Marshall



Some Members of the Research Group

Postgraduate students



Sarojinie Fernando PhD (recently submitted)



Kate Richards PhD (year 2)



Lyndal Henden Honours



Brigid Betz-Stablein PhD (year 3) Khair Jones PhD (year 1)

Staff with interest in infectious disease applications

Prof Martin Hazelton A/Prof Geoff Jones Dr Chris Jewell Dr Jonathan Marshall



Improved Spatio-Temporal Risk Estimation





- Local linear versus density ratio estimation of relative risk.
- Application to myrtle wilt in Tasmanian myrtle beech (*Nothofagus cunninghamii*).



Some Members of the Research Group

Postgraduate students



Sarojinie Fernando PhD (recently submitted)



Kate Richards PhD (year 2)



Lyndal Henden Honours



Brigid Betz-Stablein PhD (year 3) Khair Jones PhD (year 1)

Staff with interest in infectious disease applications

Prof Martin Hazelton A/Prof Geoff Jones Dr Chris Jewell Dr Jonathan Marshall



Modelling Progression of Eye Disease



- Using ideas from geographical disease mapping to examine spatial patterns of damage over retina.
- Interpretation of visual field data must account for physiology and optics.

Some Members of the Research Group

Postgraduate students



Sarojinie Fernando PhD (recently submitted)



Kate Richards PhD (year 2)



Lyndal Henden Honours



Brigid Betz-Stablein PhD (year 3) Khair Jones PhD (year 1)

Staff with interest in infectious disease applications

Prof Martin Hazelton A/Prof Geoff Jones Dr Chris Jewell Dr Jonathan Marshall



Modelling Foot and Mouth Disease in Vietnam



- Data quality issues for FMD in Vietnam.
- Appropriate level of data aggregation for modelling?
- Developing methods for identifying anomalies
 - High risk provinces.
 - Under-reporting.



Some Members of the Research Group

Postgraduate students



Sarojinie Fernando PhD (recently submitted)



Kate Richards PhD (year 2)



Lyndal Henden Honours



Brigid Betz-Stablein PhD (year 3) Khair Jones PhD (year 1)

Staff with interest in infectious disease applications

Prof Martin Hazelton A/Prof Geoff Jones Dr Chris Jewell Dr Jonathan Marshall



Bayesian Constrained Smoothing Problems



In some problems want flexible smooth fit under shape constraints.

Application to toxoplasmosis data from 34 cities in El Salvador.



Some Members of the Research Group

Postgraduate students



Sarojinie Fernando PhD (recently submitted)



Kate Richards PhD (year 2)



Lyndal Henden Honours



Brigid Betz-Stablein PhD (year 3) Khair Jones PhD (year 1)

Staff with interest in infectious disease applications

Prof Martin Hazelton A/Prof Geoff Jones Dr Chris Jewell Dr Jonathan Marshall



Methods of Simulation Based Inference



ABC Effective Population Size Estimates

- Modern complex stochastic models often have intractable distribution theory.
- Standard methods of inference infeasible.
- We are working on improving simulation based approaches (e.g. ABC).
- Application is estimation of effective population sizes in Bali using coalescent model.



TO BE CONTINUED ...

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 → 의 < ⊙