

Statistical Methods for real-time monitoring of health outcomes

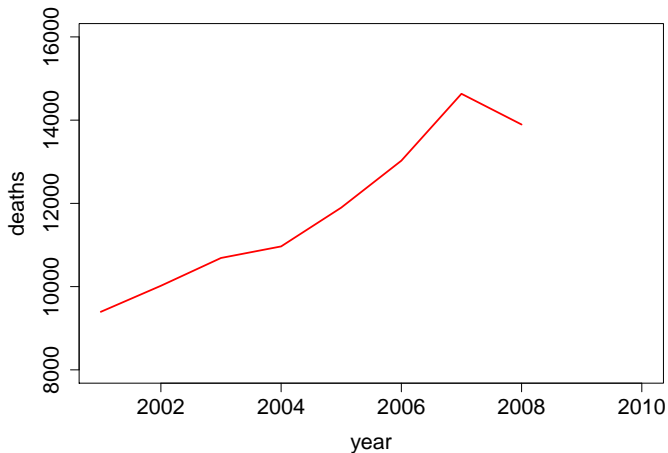
Peter J Diggle

**CHICAS, Faculty of Health and Medicine, Lancaster
University**

September 2015

- **increasing availability of electronically recorded health outcome data**
- **community and/or individual level**
- **accruing in “real-time”**
- **often spatially referenced**
- **prediction and/or explanation**
- **case-studies:**
 - **monitoring progression towards end-stage renal failure**
 - **human and veterinary surveillance of gastro-enteric illness**
 - **local-scale malaria prevalence mapping**

Chronic renal failure: UK mortality data



<http://www.endoflifecare-intelligence.org.uk>



Diagnosis, treatment and survival

Diagnosis

- Serum creatinine \Rightarrow estimated glomerular filtration rate

$$\text{eGFR} = 186 \times \left(\frac{\text{SCr}}{88.4} \right)^{-1.154} \times \text{age}^{-0.203} (\times 0.742 \text{ if female})$$

- progression can be asymptomatic for many years
- **SCr** easy to measure from blood-sample

Treatment and survival

- aggressive control of blood-pressure
- renal replacement therapy: dialysis and transplantation
- early diagnosis can slow rate of progression

	Survival rate (%) to year			
	1	2	5	10
Dialysis	79.3	64.7	33.6	10.2
Transplant (living)	98.4	96.5	90.0	76.0

Clinical guideline

Loss of $> 5\%$ eGFR per year \Rightarrow refer to secondary care

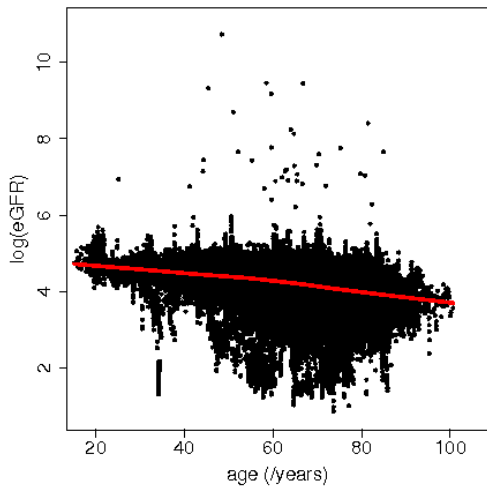
Data

- **measurements** $Y_{ij} = \log \text{eGFR}$ at **times** t_{ij} ,
explanatory variables x_i (age, sex)
 - $i = 1, \dots, m = 22,910$ “at-risk” primary care patients
 - $j = 1, \dots, n_i \leq 305$ (median $n_i = 12$)
 - $0 \leq 10.02$ years follow-up (median 4.46)
- $\mathcal{H}_i(t) = \{x_i, (t_{ij}, y_{ij}) : t_{ij} \leq t\}$

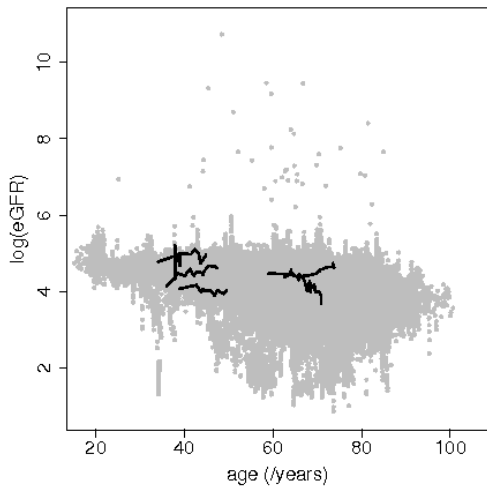
Statistical objective

$$P \left(\frac{d}{dt} \log \text{GFR} < -0.05 | \mathcal{H}_i(t) \right) = ?$$

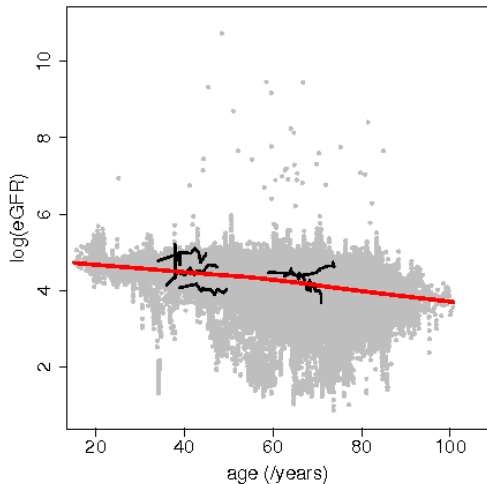
Data: all cross-sectional



Data: all cross-sectional and selected longitudinal



Data: all cross-sectional and selected longitudinal



Dynamic regression model

- **subjects** $i = 1, \dots, n$ observed at times $t_{ij}, j = 1, \dots, n_i$

$$Y_{ij} = \log(\text{eGFR})$$

- **expected value** of Y_{ij} linear in initial age and time since recruitment
- **rate of progression** varies randomly:
 - **between subjects:** random effect U_i
 - **within subjects:** random effect $C_i(t_{ij})$

Dynamic Regression Model

$$\begin{aligned} Y_{ij} &= \alpha_0 + \alpha_1 \times I(\text{female}) \\ &+ \alpha_2 \times \text{age}_{i1} + \alpha_3 \times (\text{age}_{ij} - \text{age}_{i1}) + \alpha_4 \times \max(0, \text{age}_{ij} - 56.5) \\ &+ U_i + C_i(t_{ij}) + Z_{ij} \end{aligned}$$

- Z_{ij} : measurement error, $N(0, \tau^2)$
- U_i : between-subject random intercept, $N(0, \omega^2)$
- $C_i(t)$: within-subject stochastic process

Model $C_i(t)$ as **integrated Brownian motion**

$$C_i(t) = \int_0^t B_i(u) du$$

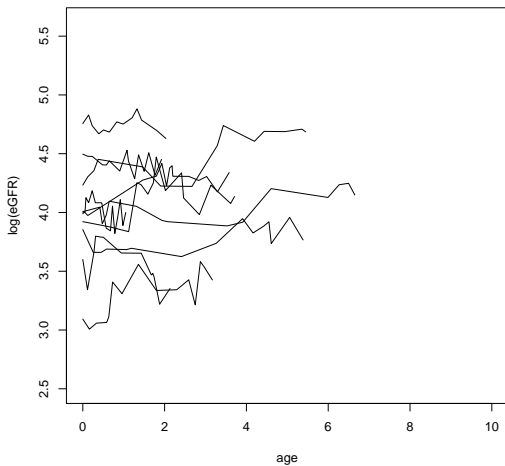
$$B_i(u) | B_i(s) \sim N(B_i(u), (u - s)\sigma^2)$$

$B_i(u)$ is rate of progression for subject i at time t

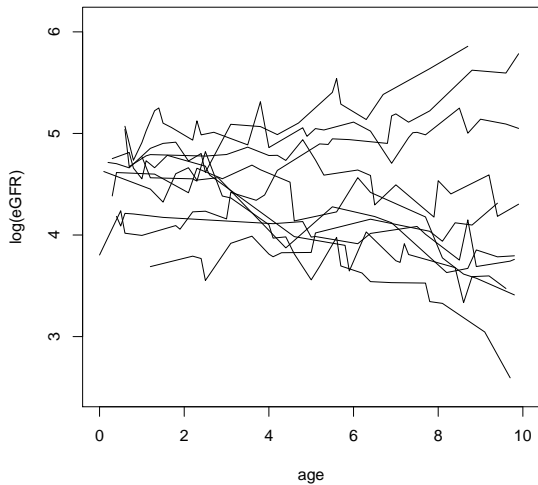
RE(%) = $100(\exp(\hat{\alpha}) - 1)$ corresponds to estimated annual percentage change in renal function.

Parameter	Estimate	SE	RE(%)
α_0 intercept	4.6006	0.0203	
α_1 female	-0.0877	0.0048	-8.4
α_2 age on entry	-0.0048	0.0004	-0.5
α_3 follow-up	-0.0232	0.0011	-2.3
α_4 age > 56.5	-0.0075	0.0006	-0.6
ω^2 intercept	0.1111	0.0012	
σ^2 signal	0.0141	0.0002	
τ^2 noise	0.0469	0.0001	

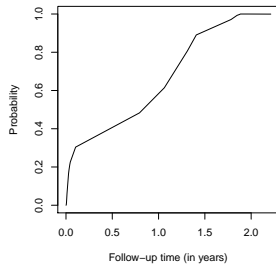
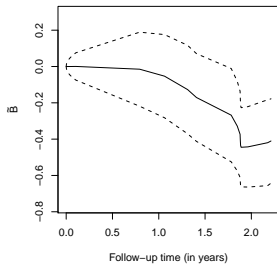
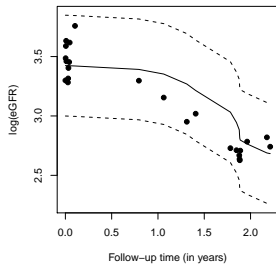
Sample data-sequences



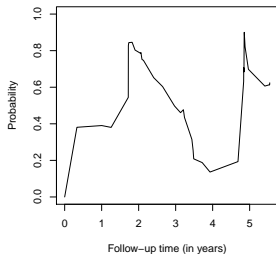
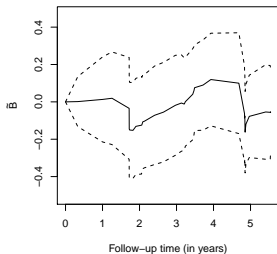
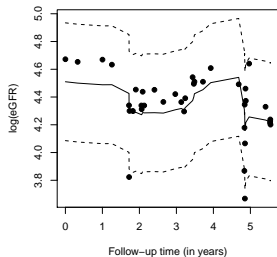
Simulations



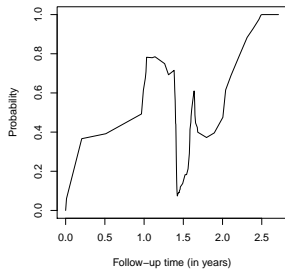
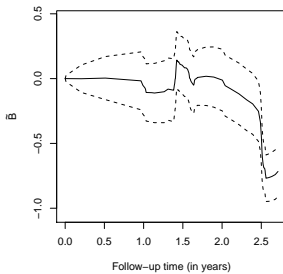
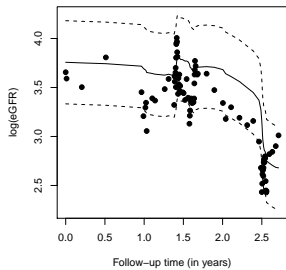
Prediction: classic progression pattern



Prediction: AKI (Acute Kidney Injury) recovery



Prediction: non-recovery from AKI



Field-testing: comparative evaluation against current methods

- eye-balling
- OLS fit to three most recent values

Informative follow-up: eGFR more likely to be measured when subject is in poor health

⇒ joint modelling of eGFR measurements and follow-up times

Implementation: in clinical practice...needs informatics expertise

Reported UK annual incidence

Campylobacter	50,000
Salmonella	10,000
Cryptosporidium	5,000
Giardia	3,000
E Coli,...	?

AEGISS: Ascertainment and **E**nhancement of **G**astroenteric **I**nfection **S**urveillance **S**tatistics

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

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

$$\begin{aligned}\text{actual} &= \text{expected} \times \text{unexpected} \\ \lambda(\mathbf{x}, t) &= \lambda_0(\mathbf{x}, t) \times R(\mathbf{x}, t)\end{aligned}$$

Scientific objective

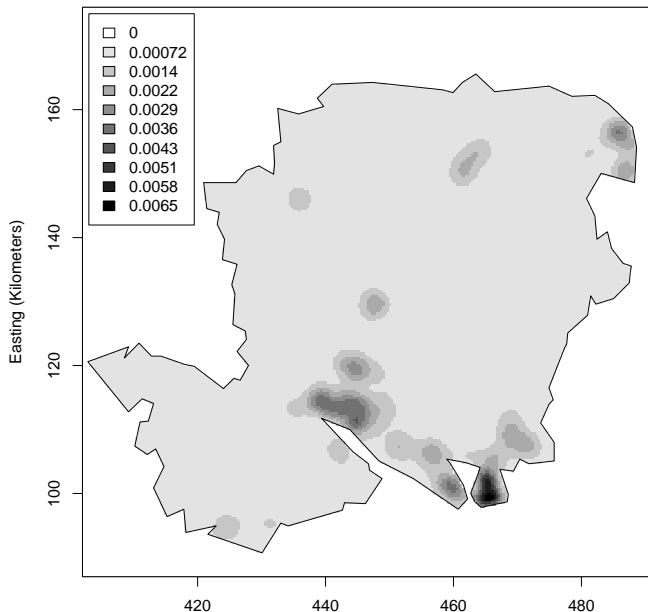
- use incident data up to time t to construct predictive distribution for current “risk” surface, $R(\mathbf{x}, t)$,
- hence identify **anomalies**, for further investigation.

$$\lambda(\mathbf{x}, t) = \lambda_0(\mathbf{x}, t)R(\mathbf{x}, t)$$

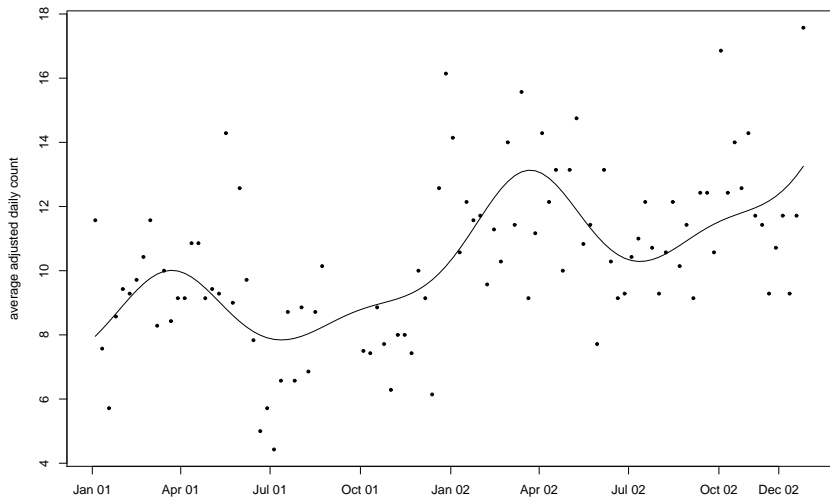
- $\lambda_0(\mathbf{x}, t) = \lambda_0(\mathbf{x})\mu_0(t)$
- $R(\mathbf{x}, t) = \exp\{S(\mathbf{x}, t)\}$
- $S(\mathbf{x}, t) =$ spatio-temporal Gaussian process

Conditional on $R(\mathbf{x}, t)$, incident cases form an inhomogeneous Poisson process with intensity $\lambda(\mathbf{x}, t)$

$\hat{\lambda}_0(\mathbf{x})$: adaptive kernel smoothing

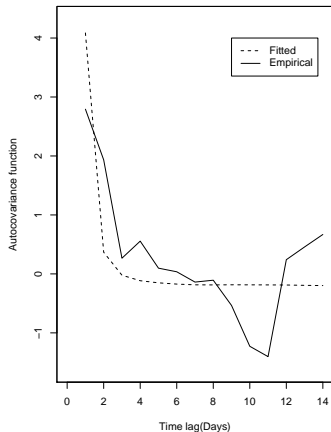
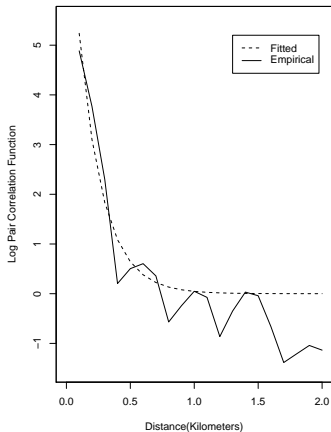


$\hat{\mu}_0(\mathbf{t})$: Poisson log-linear model



Spatio-temporal covariance

$$\rho(\mathbf{u}, \mathbf{v}) = \rho_x(\mathbf{u})\rho_t(\mathbf{v})$$

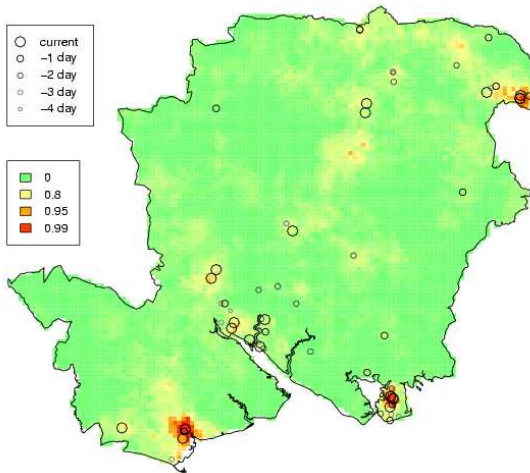


- **plug-in for estimated model parameters**
- **MCMC to generate samples from conditional distribution of $S(x, t)$ given data up to time t**
- **choose critical threshold value $c > 1$**
- **map empirical exceedance probabilities,**

$$p_t(x) = P(\exp\{S(x, t)\} > c | \text{data})$$

- **web-based reporting with daily updates**
(www.lancs.ac.uk/staff/diggle/aegiss/)

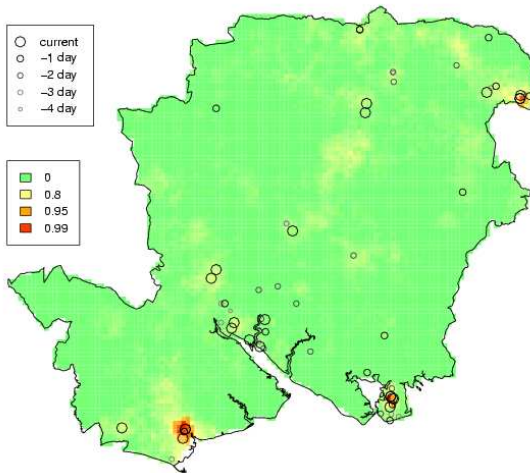
Spatial prediction: 6 March 2003



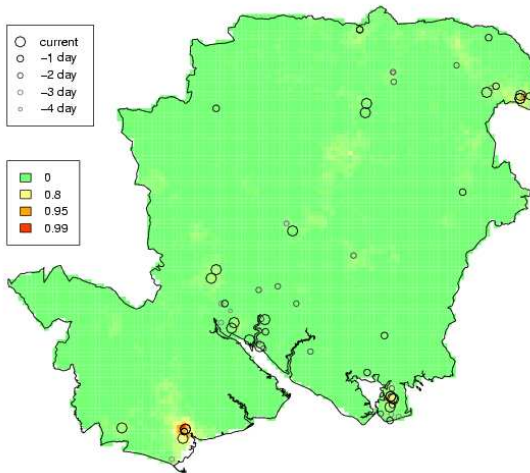
$c = 2$

Spatial prediction: 6 March 2003

$c = 4$



Spatial prediction: 6 March 2003



$c = 8$

- expand to national coverage
- integrate human and small-animal veterinary surveillance

BUT...

- replacement of single NHS Direct by multiple NHS111 services
- full post-code data no longer available!

SAVSNET: real-time data-feed from network of small-animal vet practices:

- practice location
- species (cat or dog)
- diagnosis

<http://www.savsnet.co.uk/realtimedata/>

- re-calibration of AEGISS model
- coarser spatial resolution...fitting spatially continuous models to spatially discrete data
- joint modelling of human and animal incidence
- implementation as part of routine surveillance systems

Malaria prevalence mapping



Single prevalence survey

Sample n individuals, observe Y positives

$$Y \sim \text{Bin}(n, p)$$

Multiple prevalence surveys

Sample n_i individuals, observe Y_i positives, $i = 1, \dots, m$

$$Y_i \sim \text{Bin}(n_i, p_i) ?$$

Extra-binomial variation

Sample n_i individuals, observe Y_i positives, $i = 1, \dots, m$

$$Y_i | d_i, U_i \sim \text{Bin}(n_i, p_i) \quad \log\{p_i/(1 - p_i)\} = d_i' \beta + U_i$$

Question: What to do if the d_i and/or the U_i are spatially structured

- **Latent spatially correlated process**

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

- **Latent spatially independent random effects**

$$\mathbf{U}_i \sim \text{iidN}(\mathbf{0}, \nu^2)$$

- **Linear predictor (regression model)**

$\mathbf{d}(\mathbf{x})$ = environmental variables at location \mathbf{x}

$$\eta(\mathbf{x}_i) = \mathbf{d}(\mathbf{x}_i)' \boldsymbol{\beta} + \mathbf{S}(\mathbf{x}_i) + \mathbf{U}_i$$

$$p(\mathbf{x}_i) = \log[\eta(\mathbf{x}_i) / \{1 - \eta(\mathbf{x}_i)\}]$$

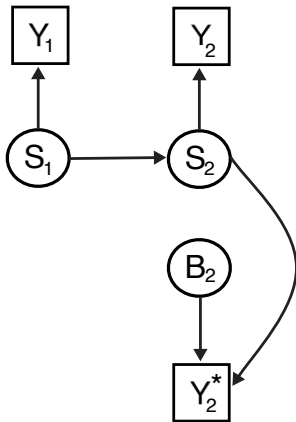
- **Conditional distribution for positive proportion \mathbf{Y}_i/n_i**

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

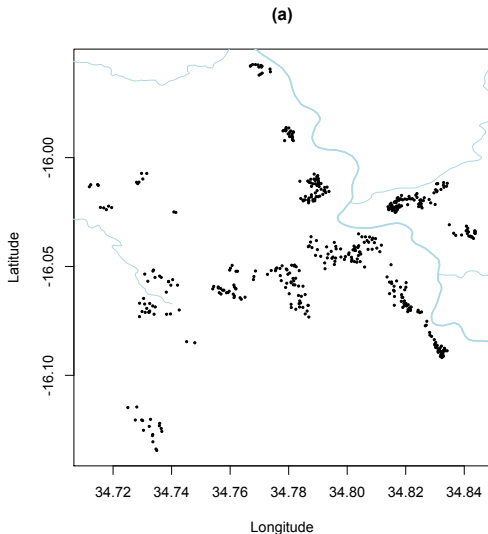
Multiple surveys (Giorgi et al, 2015)

Surveys: $i = 1, \dots, r$ **locations** $x_{ij} : j = 1, \dots, n_i$

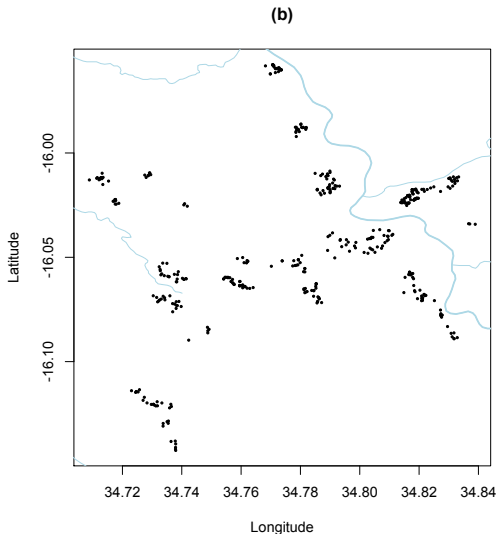
$$\eta_{ij} = \mathbf{d}(x_{ij})^\top \beta_1 + \mathbf{S}_i(x_{ij}) + \mathbf{I}(i \in \mathcal{B})[\mathbf{B}_i(x_{ij}) + \mathbf{d}(x_{ij})' \beta_i] + \mathbf{U}_{ij}$$



Malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): rMIS individual locations



Malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): eMIS individual locations



Continuous time: rolling malaria indicator surveys

Hotspots: $P(\text{prevalence} > 20\%)$

Continuous time: rolling malaria indicator surveys

Coldspots: $P(\text{prevalence} < 5\%)$

- **Operational issues**

- predictive probability of exceedance over intervention threshold
- to inform, but not to over-ride, clinical judgement

- **Methodological issues:**

- observational studies vs trials
- long series with irregular follow-up times
- informative follow-up...marked point process models

Monitoring renal function

Diggle, P.J., Sousa, I. and Asar, Ö. (2014). Real-time monitoring of progression towards renal failure in primary care patients. (submitted)

Gastro-enteric surveillance

Diggle, P.J., Rowlingson, B. and Su, T-L. (2005). Point process methodology for on-line spatio-temporal disease surveillance. *Environmetrics*, **16**, 423–34.

Malaria prevalence mapping

Giorgi, E., Sesay, S.S., Terlouw, D.J. and Diggle, P.J. (2015). Combining data from multiple spatially referenced prevalence surveys using generalized linear geostatistical models. *Journal of the Royal Statistical Society A* **178**, 445–464.