

Model-Based Geostatistics for Prevalence Mapping in Low-Resource Settings

Peter J Diggle^{1,2} and Emanuele Giorgi¹

Lancaster University¹ and University of Liverpool²



Lancaster
Medical School



UNIVERSITY OF
LIVERPOOL

INSTITUTE OF INFECTION
AND GLOBAL HEALTH

References

Diggle, P.J., Moyeed, R.A. and Tawn, J.A. (1998). Model-based Geostatistics (with Discussion). *Applied Statistics* **47** 299–350.

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.

Giorgi, E., Sesay, S.S., Terlouw, D.J. and Diggle, P.J. (2014). Combining data from multiple spatially referenced prevalence surveys using generalized linear geostatistical models. *Journal of the Royal Statistical Society A* (to appear)

Rodrigues, A. and Diggle, P.J. (2010). A class of convolution-based models for spatio-temporal processes with non-separable covariance structure. *Scandinavian Journal of Statistics*, **37**, 553–567.

Zoure, H.G.M., Noma, M., Tekle, A.H., Amazigo, U.V., Diggle, P.J., Giorgi, E. and Remme, J.H.F. (2014). The geographic distribution of onchocerciasis in the 20 participating countries of the African Programme for Onchocerciasis Control: 2. Pre-control endemicity Levels and estimated number infected. *Parasites and Vectors*, **7**, 326

Acknowledgements

CHICAS, Lancaster: Ole Christensen, Barry Rowlingson, Michelle Stanton, Ben Taylor, Rachel Tribbick

MLW, Blantyre, Malawi Sanie Sesay, Anja Terlouw

APOC, Ouagadougou: Hans Remme, Honorat Zoure, Sam Wanji

IRI, Columbia University: Madeleine Thomson

...and many others

- introduction
- general remarks on statistical modelling
- the standard binomial geostatistical model: **Loa loa**
- low-rank approximations: **river blindness**
- combining data from multiple surveys: **malaria**
- spatially structured zero-inflation: **river blindness re-visited**
- implementation
- closing remarks

Low resource settings



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, $i = 1, \dots, m$

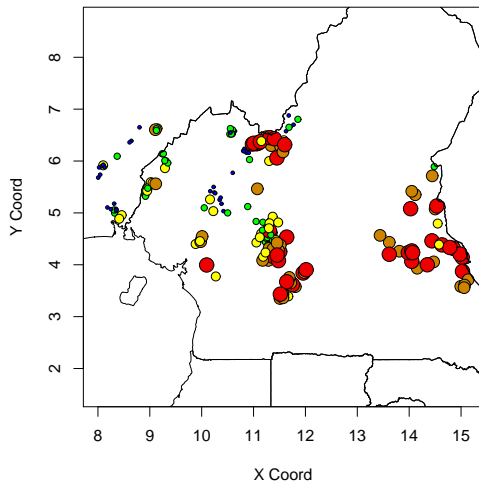
$$Y_i | d_i, U_i \sim \text{Bin}(n_i, p_i) \quad \log\{p_i/(1 - p_i)\} = d_i' \beta + U_i$$

This talk

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

- **traditionally, a self-contained methodology for spatial prediction, developed at École des Mines, Fontainebleau, France**
- **nowadays, that part of spatial statistics which is concerned with data obtained by spatially discrete sampling of a spatially continuous process**

A geostatistical data-set: Loa loa prevalence surveys



Model-based Geostatistics

(Diggle, Moyeed and Tawn, 1998)

- **the application of general principles of statistical modelling and inference to geostatistical problems**
 - formulate a model for the data
 - use likelihood-based methods of inference
 - answer the scientific question

- models are **devices to answer questions**
- models should:
 - be **not demonstrably inconsistent** with the data;
 - incorporate the underlying science, **where this is well understood**
 - **be as simple as possible**, within the above constraints

“Too many notes, Mozart”

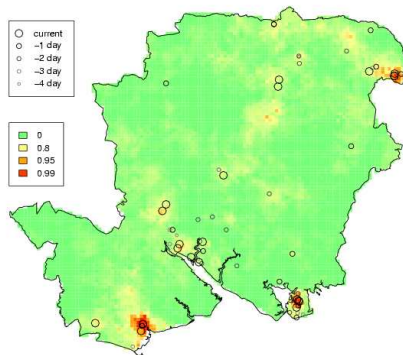
Emperor Joseph II

“Only as many as there needed to be”

Mozart (apochryphal?)

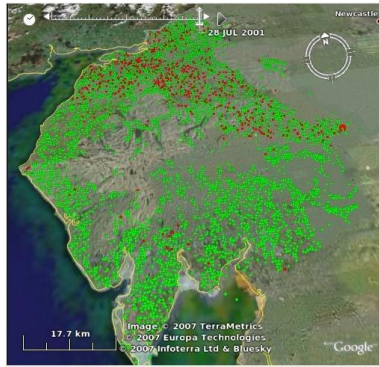
Empirical modelling: The AEGISS project (Diggle, Rowlingson and Su, 2005)

- early detection of anomalies in local incidence
- data on 3374 consecutive reports of non-specific gastro-intestinal illness
- log-Gaussian Cox process, space-time correlation $\rho(u, v)$



Mechanistic modelling: the 2001 UK FMD epidemic (Diggle, 2006)

- Predominantly a classic epidemic pattern of spread from an initial source
- Occasional apparently spontaneous outbreaks remote from prevalent cases
- $\lambda(x, t | \mathcal{H}_t)$ = conditional intensity, given history \mathcal{H}_t



Onchocerciasis (River Blindness)



A **P** **O** **C**

**frican
rogramme for
nchocerciasis
ontrol**

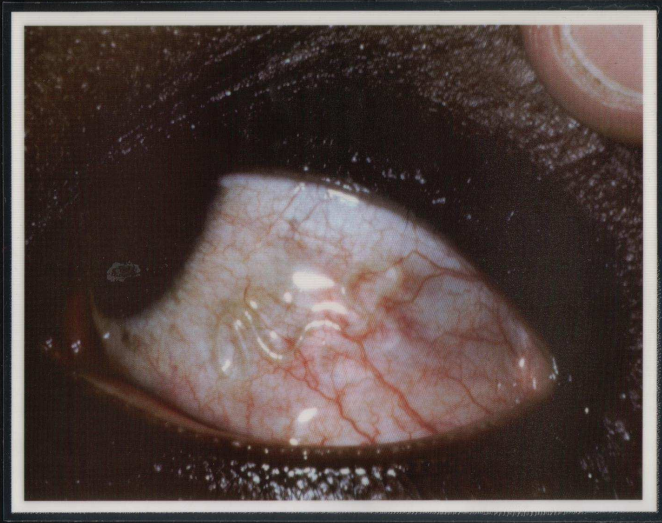
- “river blindness” – endemic in wet tropical regions
- donation programme of mass treatment with ivermectin
- approximately 60 million treatments to date, in 19 countries
- serious adverse reactions experienced by some patients highly co-infected with *Loa loa* parasites
- precautionary measures put in place before mass treatment in areas of high *Loa loa* prevalence

<http://www.who.int/pbd/blindness/onchocerciasis/en/>

Loa loa young



...and old



The Loa loa prediction problem

Ground-truth survey data

- random sample of subjects in each of a number of villages
- blood-samples test positive/negative for *Loa loa*

Environmental data (satellite images)

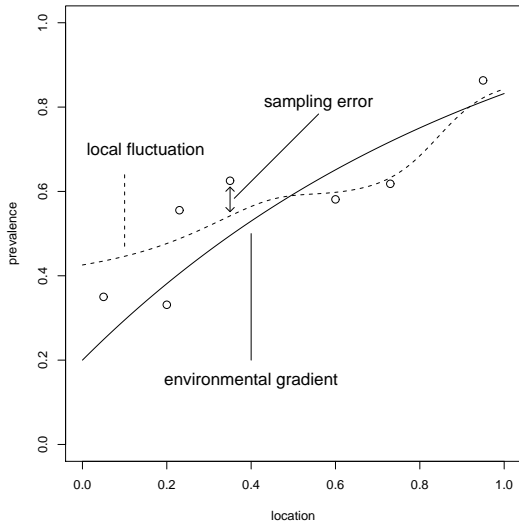
- measured on regular grid to cover region of interest
- elevation, green-ness of vegetation

Objectives

- predict local prevalence throughout study-region (Cameroon)
- compute local exceedance probabilities,

$$P(\text{prevalence} > 0.2 | \text{data})$$

Schematic representation of Loas loa model



The LoA loA modelling strategy

- use relationship between environmental variables and ground-truth prevalence to construct preliminary predictions via **logistic regression**
- use local deviations from regression model to estimate smooth **residual spatial variation**
- model-based approach acknowledges **uncertainty in predictions**

“The answer to any prediction problem is a probability distribution”

Peter McCullagh

Loa loa: a generalised linear model

- **Latent spatially correlated process**

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

- **Linear predictor (regression model)**

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

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

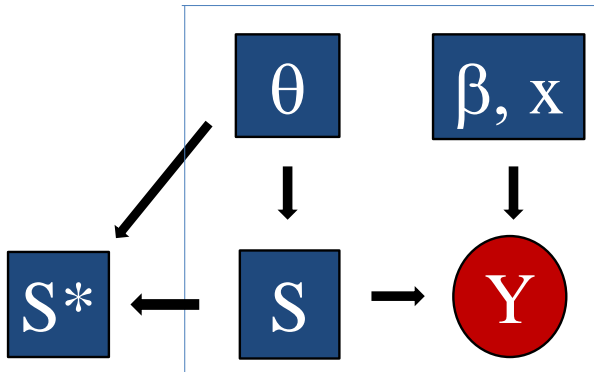
$$\mathbf{Y}_i | \mathbf{S}(\cdot) \sim \text{Bin}\{n_i, \mathbf{p}(\mathbf{x}_i)\} \text{ (binomial sampling)}$$

Conditional dependence structure

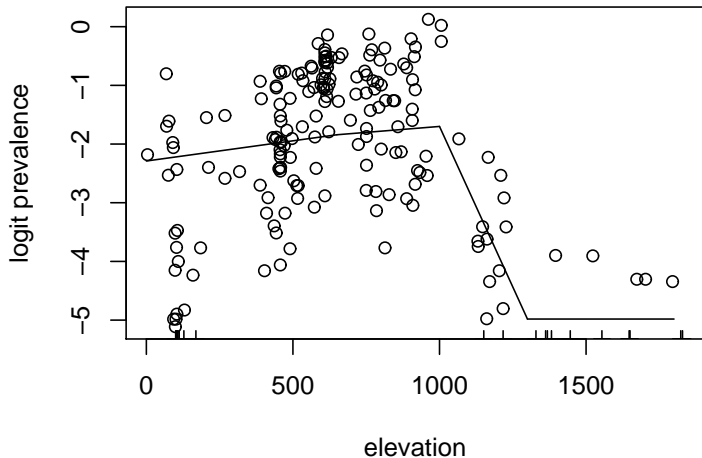
Signal: S, S^* (data-locations and prediction locations)

Data: Y (data-locations only)

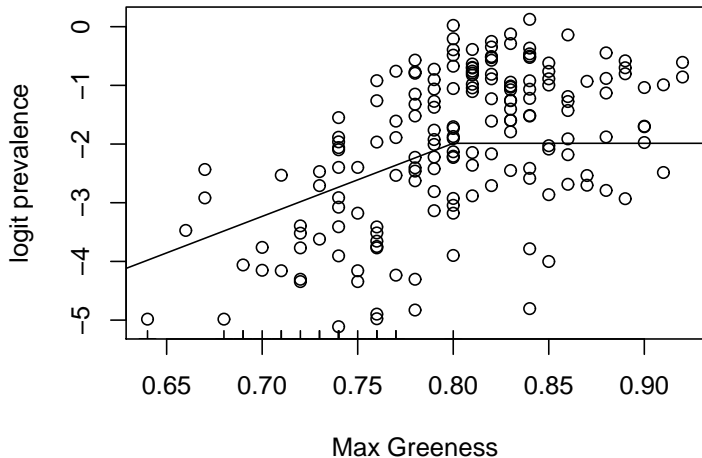
Parameters: β (regression terms), θ (covariance structure)



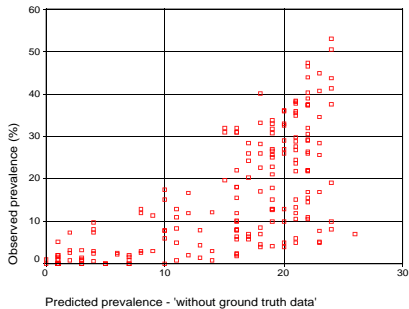
logit prevalence vs elevation



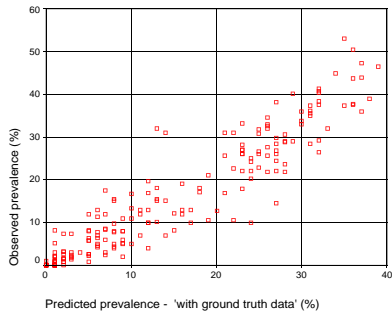
logit prevalence vs max NDVI



How useful is the geostatistical modelling?



Logistic regression



Model-based geostatistics

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

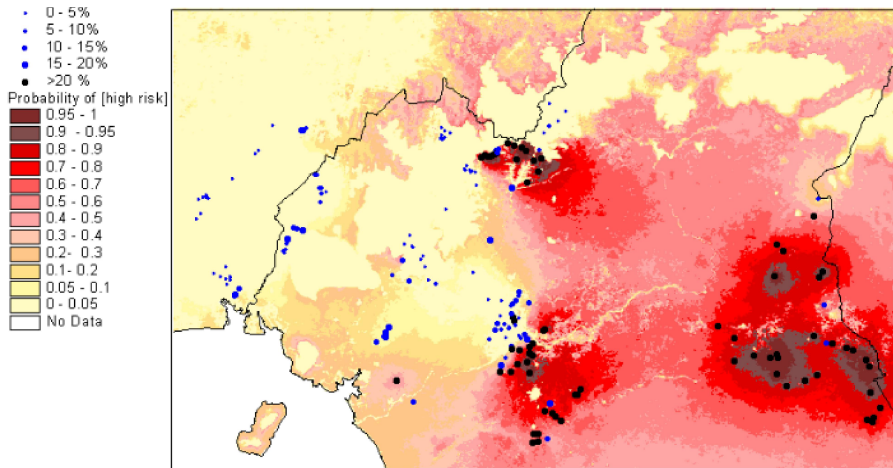


Figure 6: 'PCM for [high risk] in Cameroon based on 'ERM with ground truth data.

Extending the model

- **non-spatial extra-binomial variation**
- **low-rank approximations;**
- **combining data from multiple surveys**
 - **randomised and non-randomised**
 - **at different times**
- **spatially structured zero-inflation.**

Non-spatial extra-binomial variation

- **Latent spatially correlated process**

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

- **Latent spatially independent random effects**

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

- **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**

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

Low-rank approximations

(Rodrigues and Diggle, 2010)

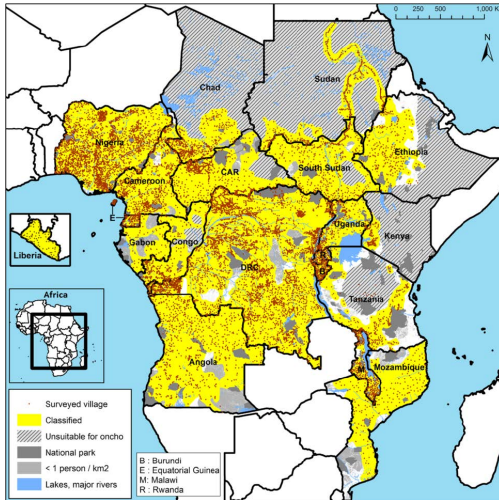
$$S(\mathbf{x}) \approx \mu + \sum_{j=1}^M w(\mathbf{x} - \mathbf{k}_j) Z_j$$

- $w(\mathbf{u})$: kernel function
- $Z_j \sim \text{iid } N(0, \nu^2)$
- $\mathbf{k}_j \in \mathbf{A} \subset \mathbb{R}^2$: fixed set of points

Choose $w(\cdot)$ to approximate to preferred family of correlation functions

Computation linear in number of prediction points

Application: onchocerciasis mapping Africa-wide (Zoure et al, 2014): 14,473 survey locations



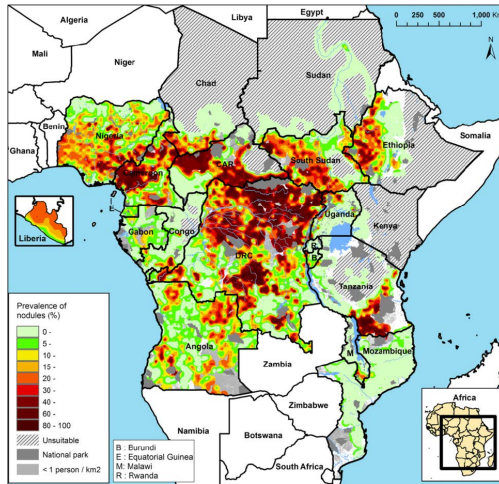
Application: onchocerciasis mapping Africa-wide (Zoure et al, 2014): low-rank model

- $M = 10734$ points X_j in regular lattice at spacing 0.1 degrees
- to approximate Matérn correlation, $M(\phi, \kappa)$, $\kappa = 2$

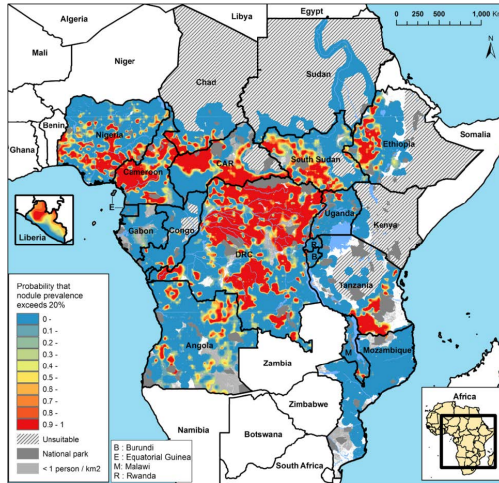
$$w(u) = \phi^{-1} \exp(-2\sqrt{2} u/\phi)$$

Parameter	estimate	95% confidence interval
μ	2:451	(2.469, 2.432)
ν^2	31:570	(31.038, 32.112)
ϕ	65:208	(64.993, 66.301)

Application: onchocerciasis mapping Africa-wide (Zoure et al, 2014): prevalence estimates



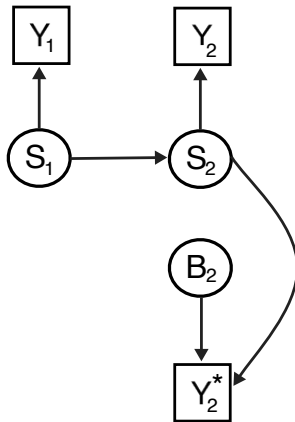
Application: onchocerciasis mapping Africa-wide (Zoure et al, 2014): exceedance probabilities



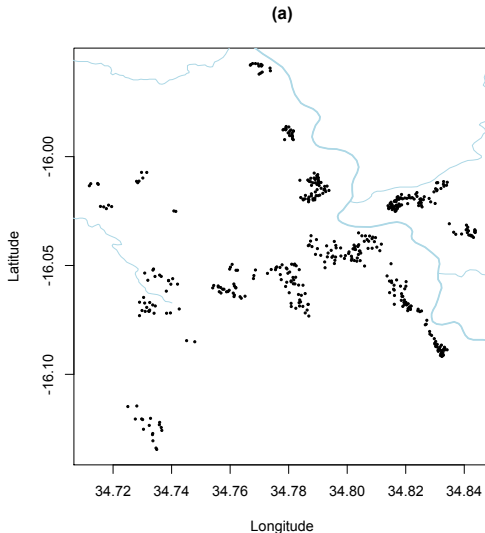
Multiple surveys (Giorgi et al, 2015)

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

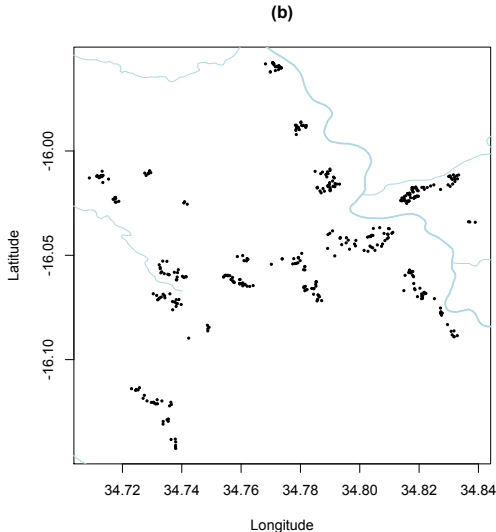
$$\eta_{ij} = \mathbf{d}(x_{ij})^\top \beta_1 + \mathbf{S}_i(x_{ij}) + \mathbf{I}(i \in \mathcal{B})[\mathbf{B}_i(x_{ij}) + \mathbf{d}(x_{ij})' \beta_i] + \mathbf{U}_{ij}$$



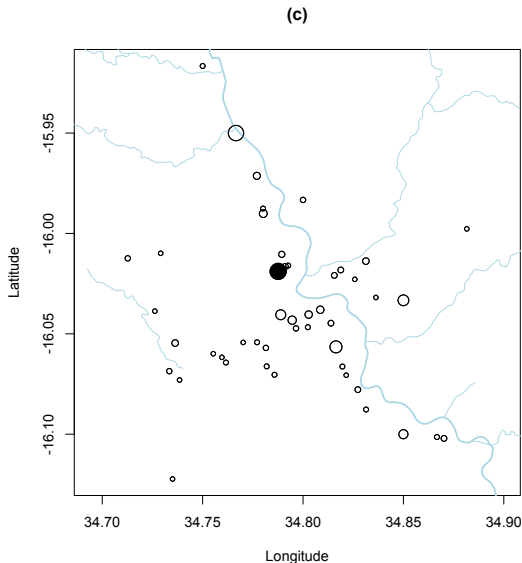
Application: malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): rMIS individual locations



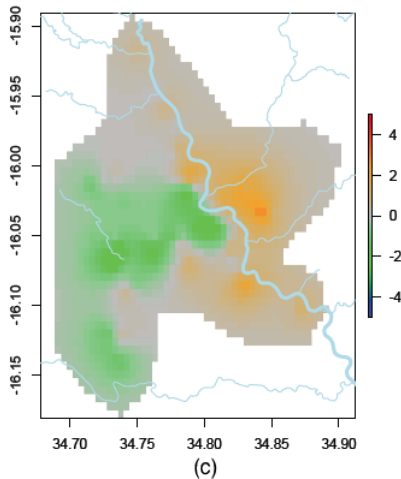
Application: malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): eMIS individual locations



Application: malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): EAG village locations and prevalences

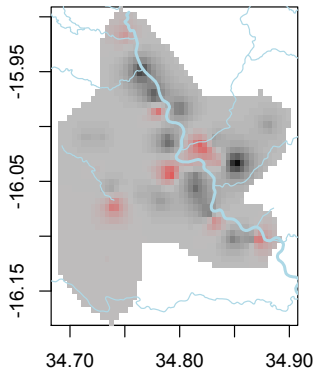


Application: malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): estimated prevalence map



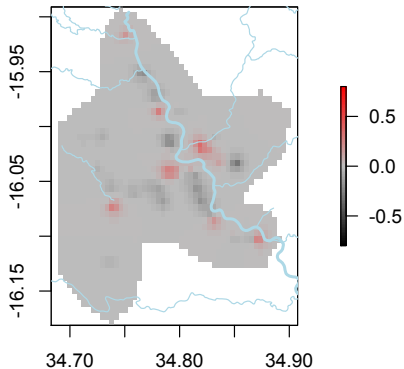
Application: malaria mapping, Chikhwawa district, Malawi (Giorgi et al, 2015): estimated bias maps

Without adjustment for SES



(a)

With adjustment for SES



(b)

Continuous time: rolling malaria indicator surveys

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

Continuous time: rolling malaria indicator surveys

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

Spatially structured zero-inflation: river blindness re-visited

- public health experts have strong sense that some areas are fundamentally unsuitable for onchocerciasis transmission
- hence need to incorporate mix of structural and chance zeros

Non-spatial model

$$Y_i \sim \begin{cases} 0 & : \text{wp } q_i \\ \text{Bin}(n_i, p_i) & : \text{wp } 1 - q_i \end{cases}$$

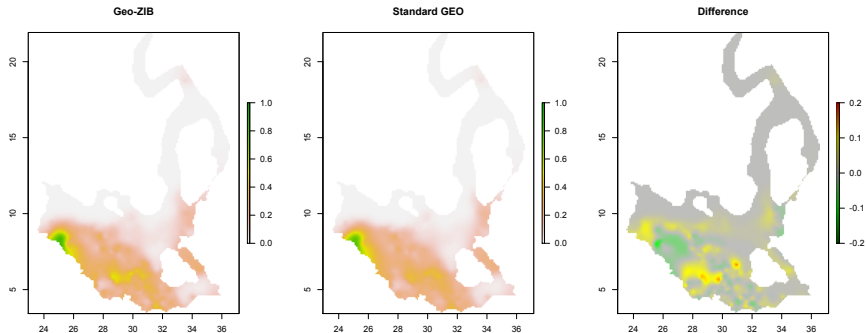
Spatial model

$\{q_i, p_i\} \rightarrow \{Q(x), P(x)\} : x \in \mathbb{R}^2 \sim \text{bivariate stochastic process}$

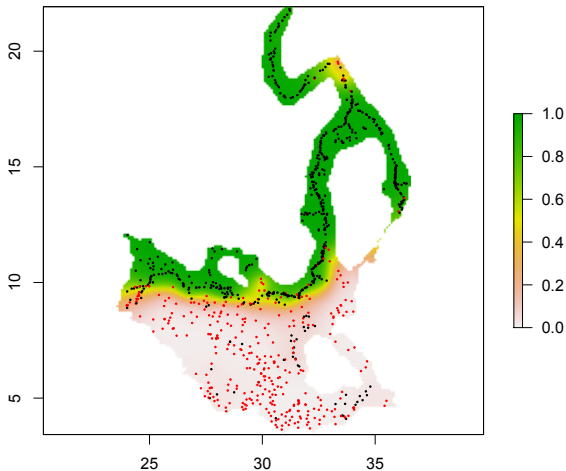
$$P(Y = y | S_1(x), S_2(x)) = \begin{cases} Q(x) + (1 - Q(x)) \times \text{Bin}(0; n, p(x)) & : y = 0 \\ (1 - Q(x)) \times \text{Bin}(y; n, p(x)) & : y > 0 \end{cases}$$

- $\text{logit}(Q(x)) = \mu_1 + S_1(x)$
- $\text{logit}(P(x)) = \mu_2 + S_2(x)$
- $\{S_1(x), S_2(x)\} \sim \text{bivariate Gaussian process}$

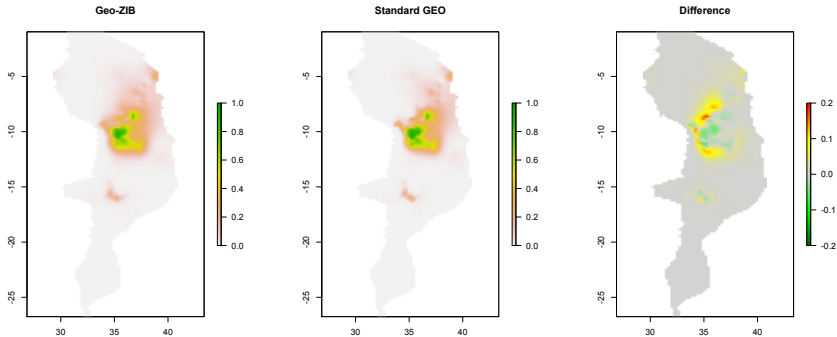
Sudan: probability exceedance map



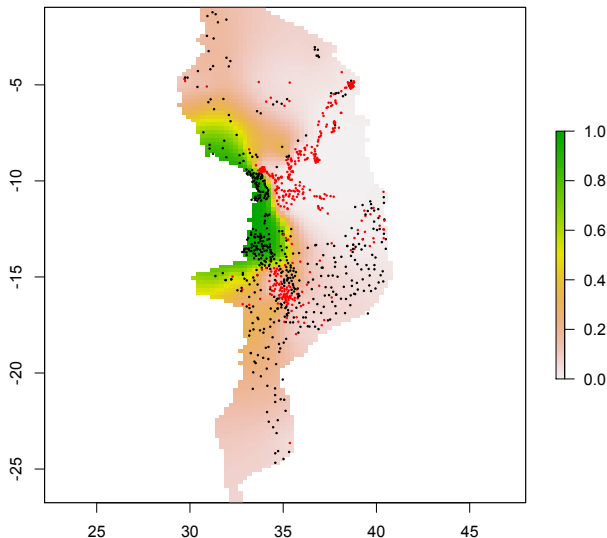
Sudan: non-transmissible probability map ($Q(x)$)



Mozambique/Malawi/Tanzania: probability exceedance map



Mozambique/Malawi/Tanzania: non-transmissible probability map ($Q(\mathbf{x})$)



- Monte Carlo maximum likelihood
- Plug-in prediction
- R package `PrevMap`
- Bayesian version?

Closing remarks

- **principled statistical methods**
 - make assumptions explicit
 - deliver optimal estimation within the declared model
 - make proper allowance for predictive uncertainty

- but there is no such thing as a free lunch

“We buy information with assumptions”

C H Coombs

- which is why statistics is at its most effective when conducted as a **dialogue with substantive science**
- and this should **guide the way we teach statistics** ...especially to science students