

# Achieving explanatory depth and spatial breadth in infectious disease modelling: Integrating active and passive case surveillance

Statistical Methods in Medical Research 0(0) 1–15 © The Author(s) 2019 Article reuse guidelines: sagepub.com/journals-permissions DOI: 10.1177/0962280219856380 journals.sagepub.com/home/smm

**SAGE** 

# Luca Nelli, D Heather M Ferguson and Jason Matthiopoulos

## Abstract

Ideally, the data used for robust spatial prediction of disease distribution should be both high-resolution and spatially expansive. However, such in-depth and geographically broad data are rarely available in practice. Instead, researchers usually acquire either detailed epidemiological data with high resolution at a small number of active sampling sites, or more broad-ranging but less precise data from passive case surveillance. We propose a novel inferential framework, capable of simultaneously drawing insights from both passive and active data types. We developed a Bayesian latent point process approach, combining active data collection in a limited set of points, where in-depth covariates are measured, with passive case detection, where error-prone, large-scale disease data are accompanied only by coarse or remotelysensed covariate layers. Using the example of malaria, we tested our method's efficiency under several hypothetical scenarios of reported incidence in different combinations of imperfect detection and spatial complexity of the environmental variables. We provide a simple solution to a widespread problem in spatial epidemiology, combining latent process modelling and spatially autoregressive modelling. By using active sampling and passive case detection in a complementary way, we achieved the best-of-both-worlds, in effect, a formal calibration of spatially extensive, errorprone data by localised, high-quality data.

#### Keywords

Bayesian modelling, disease mapping, imperfect detection, latent point process, N-mixture models, spatial epidemiology

## I Introduction

Predictive maps of disease risk, typically obtained by modelling the spatial heterogeneity in disease incidence as a function of underlying covariates, can be crucial for targeting effective control and surveillance.<sup>1–6</sup> However, reliable prediction at the landscape scale is often hindered by lack of appropriate, high-resolution spatial data. Traditionally, incidence data and potential explanatory covariates are collected either systematically – using active sampling by researchers – or opportunistically – from clinical records reported at health facilities. Each of these sampling strategies has its own limitations.<sup>7</sup> For example, by collecting detailed data for both disease incidence and related covariates, data from active sampling allows models to achieve high explanatory power but not to make large-scale extrapolation and predictions in areas where fine scale covariates are not directly measurable.<sup>8,9</sup> On the other hand, passive sampling yields data from a large number of geographically dispersed cases which are more amenable for large-scale predictions, but these data often suffer from severe reporting biases<sup>10–13</sup> and can be paired with only coarse environmental covariates that have limited explanatory power.<sup>4</sup> As the drawbacks of one strategy are clearly the strengths of the other, modelling frameworks that consider these two types of data simultaneously and complementarily would strengthen our biological insights and predictive power.

Institute of Biodiversity Animal Health and Comparative Medicine, University of Glasgow, Glasgow, UK

**Corresponding author:** Luca Nelli, Institute of Biodiversity Animal Health and Comparative Medicine, University of Glasgow, University Avenue, Graham Kerr Building, Glasgow G12 8QQ, UK. Email: luca.nelli@glasgow.ac.uk Active sampling is typically conducted by research teams that focus on a small number of predetermined locations, with collection of detailed environmental or epidemiological variables including clinical samples,<sup>9,14–19</sup> entomological indicators (for vector-borne disease),<sup>17,18,20–22</sup> human demographic and socioeconomic factors<sup>9,19,23,24</sup> or fine-scale environmental conditions.<sup>25,26</sup> Such data can provide high power for explaining variation in risk across focal sites,<sup>7</sup> but lack predictive breadth across space because many of the crucial covariates are not available for unsampled locations.<sup>9</sup>

Clinical records from passive case detection offer the potential of expansive descriptions of spatial incidence patterns. However, since these incidence data often arise from self-reporting at health centres, they can be biased by their opportunistic nature. Reporting bias is well acknowledged for numerous infectious disease systems<sup>27–29</sup> and can be expressed as a combined function of distance from health facilities, the likelihood of asymptomatic cases and sociodemographic factors<sup>10–13,30–35</sup> or more complex measures of travel time.<sup>36</sup> Despite this limitation, health centre surveys remain the primary source of information for disease monitoring. Another drawback of spatial models of incidence data gathered from passive case detection, relates to the availability of environmental predictor data. If the locality of the patient is recorded, incidence data can be spatially plotted but researchers and public health workers are unlikely to be able to directly measure some detailed explanatory variables at those localities. Therefore, when modelling the incidence data, only large-scale but coarse layers are customarily considered. While these bring more geographically expansive information than the highly localised survey data, they generally consist of remotely sensed covariates and summary records such as bioclimatic, geomorphological, vegetation indexes, human population density or road networks,<sup>9,37–40</sup> that typically contribute limited explanatory power.

Some studies<sup>7,14–16,23,41–44</sup> make use of data from both active and passive case detections together, but focus on independent analysis and comparison of results from these separate data sources rather than integrating them. Analysing these two data sources jointly can be viewed as challenging<sup>15</sup> because their limitations imply a trade-off between explanatory depth and predictive breadth. However, there is clearly an opportunity to achieve complementarity by analysing them on an integrated inferential framework. Here, for the first time, we develop a spatial statistical model combining these two sources of incidence data to harness the maximum amount of information for explanatory and predictive objectives.

Our framework takes a novel approach to both the response and the explanatory variables. The dual nature of the incidence data requires specification of a statistical model that considers two different aspects of likelihood, one for the localised but precise survey data, and another for the spatially extensive but imperfect clinical reporting data. We build this part of the approach on two cornerstones of the statistical literature: the point process model<sup>45,46</sup> and the methodology of point transects.<sup>47</sup> Point processes model events (e.g. infection cases) that occur continuously in space according to an unknown intensity (a spatial surface to be estimated as a function of covariates). We observe these events as arising from two different point transects, each having its own spatially heterogeneous observation model. The first type of observation point is the active sampling location, where cases are detected near-perfectly, but only for that particular set of geographical coordinates. The second type of observation point is a clinic, where cases of the disease are reported from a broad geographical region but with probabilities of detection that decay with distance from the clinic. Regarding the explanatory variables, some environmental variables are easily collectable for both passive case detection and active sampling points, but more important and powerful variables may be available only at the latter. By importing ideas from latent process modelling<sup>48,49</sup> we use the spatially extensive clinical data together with the data-rich survey data to reconstruct latent covariates that may be hidden from direct or remote observation.

To validate the ability of our model to retrieve correct parameter values, we require these scenarios to be accompanied by known intensity surfaces for both incidence and latent explanatory variables. These requirements cannot be satisfied by real data sets, so here we have acquired our scenarios via realistic simulation, motivating our examples from a real system of a vector-borne disease. To illustrate the generality of our approach, we have hypothesised multiple contrasting scenarios of reporting bias and spatial distribution of the latent process underlying disease incidence.

We chose malaria in West Africa as an ideal example of an important environmentally-dependent infectious disease, <sup>50,51</sup> for which human exposure and infection risk is highly spatially heterogeneous and dependent on crucial environmental variables that influence interactions between people, mosquitoes, and parasites.<sup>40,52</sup> Control measures such as long-lasting insecticide treated nets (LLINs) have been crucial for impeding contact between mosquitoes and people, and have led to substantial declines in malaria prevalence across Africa in the last decade.<sup>50,51</sup> However, the success of such an approach may be undermined by development of insecticide resistance in mosquitoes, particularly in West Africa where rates are among the highest in the world.<sup>53–55</sup>

Copious and widespread data on reported cases are often available from clinics (see for example www. malariasurveys.org or www.dhsprogram.com), but detailed information on mosquito vector ecology and insecticide resistance is only available for a limited number of sites.<sup>54,56–58</sup> These challenges exist for many other vector-borne diseases whose transmission is dependent on an ecological reservoir and rely on insecticide use for control, such as for example dengue, Zika and chikunguya viruses,<sup>59</sup> Lyme and other tick-borne disease,<sup>60</sup> schistosomiasis,<sup>61</sup> Rift Valley Fever,<sup>62</sup> human African trypanosomias<sup>63</sup> or West Nile Virus.<sup>64</sup>

# 2 Methods

# 2.1 Modelling approaches

For a given area of interest subdivided into a regular grid, we consider as our sampling unit the grid cell  $i \in \{1, ..., K\}$  We first assume an underlying stochastic process f that generates numbers of cases  $N_i$  according to an underlying, spatially heterogeneous rate  $\lambda_i$ . We also assume an observation process g that allows a subset of the  $N_i$  cases to be reported at different sampling stations. We distinguish between two types of sampling stations: S is the number of active sampling points (about which we are assuming a perfect and exclusive detection but at a small distance, i.e. within the cell that contains them). We denote by J the number of clinics (about which we are assuming an imperfect but long-ranging detection). The observation process g is therefore generating the vector of incidence data reported in each  $i^{th}$  cell at different stations  $I_i = \{I_{1,i}, \ldots, I_{S,i}, I_{(S+1),i}, \ldots, I_{(S+J),i}, U_i\}$ , given the vector of probabilities  $P_i = \{P_{1,i}, \ldots, P_{S,i}, P_{(S+1),i}, \ldots, P_{(S+J),i}, Q_i\}$ .  $U_i$  represents the number of completely unreported cases in each  $i^{th}$  cell (which is a missing value in the data), given the probability  $Q_i$  of not reporting.

The general likelihood function of our models can be expressed as follows

$$L = \prod_{i=1}^{K} f(N_i | \lambda_i) g(\boldsymbol{I}_i | \boldsymbol{P}_i, N_i)$$
(1)

We built our approach incrementally, developing three distinct modelling approaches with an increasing level of complexity to allow comparison between the routes that might have traditionally been followed to analyse data arising from active sampling (model 1) and passive case detection (model 2) with our new proposed route (model 3), which reconstructs the latent processes and estimates the emergent patterns of disease incidence with increased precision and accuracy.

#### 2.1.1 Model 1 – active sampling data only

Here, we consider data that would be collected from active sampling at just a limited number *S* of active survey sites. To analyse the relationship between disease incidence and detailed measures of covariates at a set of predetermined survey points, model 1 takes the form of a Poisson Generalised Linear Model without any spatially explicit component.

Although this is a straightforward model to fit using likelihood-based libraries in all statistical platforms, we fitted it using Bayesian methods<sup>65</sup> for consistency in the comparison with models 2 and 3 that follow. The response variable is the number of observed diseases cases  $N_i$  at the location of the *i*<sup>th</sup> survey. We assume here (for simplicity, but with no loss of generality) that all the cases at the survey location are recorded (hence, a local detection probability of 1 for each case), although we acknowledge that with conventional diagnostic tests some percentage of cases can be missed.<sup>66</sup> If data are available on diagnostic sensitivity and specificity, our method can be readily extended by incorporating false negatives or positives.

The model takes the form

$$N_i \sim Poisson(\lambda_i)$$
 (2)

where the rate  $(\lambda_i)$  of disease incidence is

$$\ln(\lambda_i) = \beta_0 + \sum_{k=1}^n \beta_k X_{ik}$$
(3)

The linear predictor on the right-hand side of this expression comprises a set of *n* coefficients  $\beta$  and *n* explanatory variables *X* measured at the *i*th survey location.

Equations (2) and (3) can be generalised to take better account of specific features of the data. For example, it may be relevant to use overdispersed forms of the likelihood (relaxing the Poisson assumption) or more complicated functional forms of the linear predictor, involving polynomials, interactions or splines.

## 2.12 Model 2 - passive case detection only

Here, we considered only data coming from passive case detection. This model maintained the basic structure of model 1, i.e. it is a Bayesian Poisson regression, with reported disease cases at human dwellings or communities surrounding the health centres as the response variable and the set of environmental variables as predictors. Under our scenarios, we assumed that one of the key predictor variables (insecticide resistance *IR*, see section 2.2) could only be measured experimentally in active sampling sites, therefore we could not include it in equation (3).

We introduced the estimation of bias in reporting disease cases given by the distance from the health centres, borrowing concepts from distance sampling theory,<sup>47</sup> a group of methods, widely used to estimate the absolute abundance or spatial density of animal or plant populations. The key underlying concept is the estimation of a detection function (P(d)), which represents the decay in the probability of detecting an object with increasing distance (d) from the observer. Given the detection function and encounter rate, the absolute density of a population can be modelled at a given point, assuming perfect detection at the location of the observer P(0) = 1. In our application, this has the interpretation that if a case arises in the immediate vicinity of the clinic ( $d \approx 0$ ), then it is certain to be reported. A plausible, but flexible decay function is fitted to paired data of detections and distances. For example, detection of a malaria case from the *i*<sup>th</sup> location at the *j*<sup>th</sup> clinic, can be modelled as a half-normal function of distance from the health centre  $d_{i,j}$ , by the following<sup>47</sup>

$$p(d_{i,j}) = \exp\left(-\frac{d_{i,j}^2}{2\sigma^2}\right) \tag{4}$$

where  $\sigma$  is the shape parameter of the half-normal function (regulating how quickly the detection probability drops with distance). The distance *d* can be Euclidean, or a more complicated function of accessibility (e.g. affected by proximity between points along a given road network).

Any given case may be reported to any one of the available clinics, but clinics nearby are more likely to receive the report. The probability of any one case being reported to any one clinic (accounting for other clinics) can be modelled in terms of the distances of all the clinics from the point of occurrence of the case, as follows

$$P_{i,j} = \frac{p(d_{i,j})}{\sum_{j=1}^{J} p(d_{i,j}) + Q_i}$$
(5)

The denominator here represents all possible outcomes, i.e. the probabilities that the case is reported to any one of J centres, and the probability  $P_{Q_i}$  that the case goes completely unreported:

$$Q_i = \prod_{j=1}^{J} [1 - p(d_{i,j})]$$
(6)

Note that  $P_{i,j}$  is the standardised form of  $p(d_{i,j})$ . In fact,  $p(d_{i,j})$  is the probability of a case being reported at a given clinic (considered in isolation), purely as a function of distance, whereas  $P_{i,j}$  is the probability of reporting at a clinic, accounting for the effects of other clinics that are "competing" for the same reports and including  $Q_i$ , that is the probability of a case not being reported at all.

The likelihood of a data set comprising clinic reports may then be written as a multinomial process. In particular, for a given number of actual cases  $N_i$  (see equation (2)), the likelihood of reported disease cases  $I_i$  in the *i*<sup>th</sup> cell for the *J* clinics in the dataset is determined by the detection probabilities  $P_i$  that are function of distances between the *i*<sup>th</sup> location and the clinics, by

$$\boldsymbol{I}_i \sim Multinomial(N_i, \boldsymbol{P}_i) \tag{7}$$

where  $P_i = \{P_{1,i}, \ldots, P_{J,i}, Q_i\}.$ 

Fitting model 2 to the data yielded estimates of the shape parameter of the detection function (equation (4)) and parameters of equation (3). Although it had no spatially explicit component, we used model 2 to generate a reconstruction of the patterns of incidence across space, based on the coarse-level environmental covariates.

Hence this model did not benefit from the fine-resolution covariates that could only be measured by detailed experimental methods at survey points.

#### 2.1.3 Model 3 – active and passive data combined

The process and observation model for this joint approach to data took the form of equations (2) and (7), respectively. However, just like in model 1, equation (3) used the full set of predictors, including the partlylatent variable (i.e. insecticide resistance, available only for active sampling points but not for regions of passive case detection data collection and the rest of space). Our model for the latent variable *IR* postulated a spatial autocorrelation structure,<sup>67</sup> implying that even though we may not know the values of the latent variable at two points in space, we can express a relationship about their expected degree of similarity. Any pair of *K* cells in our grid, say  $i \in \{1, ..., K\}$  and  $k \in \{1, ..., K\}$ , was assumed to have a covariance, specified as a decreasing function of their distance

$$cov_{i,k} = \exp(-\rho d_{i,k}) \tag{8}$$

with  $\rho \ge 0$ . Again, this is one of many possible structures and our overall approach is not constrained to this functional form.

The distribution of the latent variable  $IR = \{IR_1, ..., IR_K\}$  in all the K cells, was therefore modelled as a Gaussian field from an *m*-dimensional multivariate normal distribution, where each of the dimensions represented the probability density of a cell in space

$$IR_i \sim MVN(\mu, \sum)$$
 (9)

Here, the mean vector  $\mu$  has length K (the total number of cells in geographical space), and  $\sum$  is a  $K \times K$  spatial covariance matrix<sup>68</sup> with values of 1 on the diagonal and values  $cov_{i,k}$  for the *i* row and *k* column from equation (8).

Model 3, hence, is fitted exactly as model 2 according to equation (7), but the linear predictor function (equation (3)) included all the covariates, unlike model 2, which was missing the covariate of IR. In particular, IR observations were used where available (at active sampling points), assuming that they were realisations from equation (9).

# 2.2 Model validation

We used simulated data on malaria incidence and insecticide resistance within the primary mosquito vectors to validate our models. Our specific validation aims were to (1) evaluate the match between the posterior distribution of the coefficients and the simulation process that generated the data; (2) estimate bias in reporting the clinical data as a function of distance between the location of a clinic and the village where the patient resides; and (3) recreate the missing covariate of insecticide resistance  $\widehat{IR}$  and to reconstruct the true incidence  $\hat{N}$ .

Our simulation borrowed its setting from a study currently ongoing in Southwest Burkina Faso (MiRA -Malaria in Insecticide Resistant Africa, Wellcome Trust 200222/Z/15/Z). The study covers an area of approx. 6000 km<sup>2</sup> in the health district of Banfora in south-western Burkina Faso, comprising primarily West Sudanian savannah which experiences a rainy season from May to October with little rain in other months. Malaria transmission is stable throughout the year but peaks from May to November. The major vectors are Anopheles gambiae and An. funestus. Like many other areas of Africa, the primary malaria control strategy is long lasting insecticidal nets (LLINs) that are distributed at high coverage across the country (Burkina Faso National Malaria Control Program, unpublished data). In contrast to some areas of Africa, recent LLIN distribution campaigns have had little impact on malaria prevalence and it is hypothesised that this may be due to high levels of insecticide resistance in local vector populations,<sup>69</sup> which are amongst the highest on record. Resistance to pyrethroid insecticides is widespread. Mortality after exposure (defined by the World Health Organization (WHO) as the response to the stipulated discriminating dose of permethrin) ranges from 5 to 20%.<sup>20</sup> For the purposes of data simulation, we assume that active sampling of malaria infections and insecticide resistance levels is carried out in 12 villages, and that passive case data is available from patients reporting to from 8 health centres distributed throughout the study area. This number and distribution of passive and active sampling site was selected to represent the distribution of health facilities and likely maximum amount of active survey data available.

For the simulation, we considered a square grid with a  $1 \text{ km}^2$  resolution covering the study area. We generated a dataset with reported incidence in each cell of the grid under a binomial *N*-mixture model<sup>70,71</sup> by combining two different processes: a state model, i.e. the biological process that generates malaria infection cases, and an observation model, i.e. the process that affects the probability that infection cases are reported to a health centre.

To simulate the biological process, we considered the average altitude in the cell, average yearly temperature (TEMP), annual rainfall (RAIN), human density (HUM), normalised difference vegetation index (NDVI) and insecticide resistance (IR) in mosquitoes as potential predictors.<sup>9,38,72–77</sup>

Temperature and rainfall were derived from the WorldClim database (www.worldclim.org). NDVI values were obtained using the package *MODIStsp* for R.<sup>78</sup> To create the layer of human density, we used a kernel density estimation<sup>79</sup> using GPS points of the villages (307) in the study area and the population census in each village (1755±1804 mean±dev. std., Institut national de la statistique et de la démographie, *unpublished data*) as weight field. Kernel bandwidth was chosen so as to minimise the least-squares cross validation score  $(h_{lscv})$ .<sup>80</sup>

Insecticide resistance reporting has improved over time, and global maps of insecticide resistance at coarse resolutions are now becoming available.<sup>77</sup> However, little is known about its spatial distribution at local scale.<sup>81</sup> Therefore, to explore out model's ability to retrieve latent variables of differing spatial complexity, insecticide resistance was simulated by hypothesising three different scenarios of increasing spatial autocorrelation, with parameter  $\rho$  of equation (8) set, respectively, to  $\rho_1 = 3.0$ ,  $\rho_2 = 0.7$  and  $\rho_3 = 0.3$  (Figure 1, *IR1*, *IR2*, *IR3*).

The number of malaria cases, or true incidence, in each cell  $(N_i)$  was assumed to have a positive relationship with temperature, <sup>9,74</sup> rainfall, <sup>9,73,74</sup> human density, <sup>9</sup> NDVI<sup>9,37,38,72</sup> and insecticide resistance, <sup>76</sup> and was simulated from equation (2) using the linear predictor

 $\log(\lambda_i) = \beta_0 + \beta_{HUM} HUM_i + \beta_{NDVI} NDVI_i + \beta_{RAIN} RAIN_i + \beta_{TEMP} TEMP_i + \beta_{IR} IR_i$ 

We set the equation's coefficients to the values  $\beta_0 = 2.90$ ,  $\beta_{HUM} = 0.50$ ,  $\beta_{NDVI} = 0.30$ ,  $\beta_{RAIN} = 0.20$ ,  $\beta_{TEMP} = 0.25$ ,  $\beta_{IR} = 0.50$ . Having three distinct scenarios of insecticide resistance *IR*1, *IR*2 and *IR*3, we obtained three scenarios of malaria infection cases  $N1_i$ ,  $N2_i$  and  $N3_i$ .

For the observation process, we accounted for simulated bias in reporting cases in each cell of the grid, by considering a probability of reporting as a function of the distance between a given cell and each health centre. We set the detection probabilities in each cell P(i,j) in accordance with equation (4) with  $p(d_{i,j})$  being the Euclidean distance between the centroid of the *i*<sup>th</sup> cell of the grid and each *j*<sup>th</sup> health centre. We employed three different shapes of the detection function, using different values of the shape parameter  $\sigma_A = 10$ ,  $\sigma_B = 15$ ,  $\sigma_C = 20$  (Figure 1, P<sub>A</sub>, P<sub>B</sub> and P<sub>C</sub>). Probability of reporting at active sampling stations was deliberately set at 1, to ensure that all the infection cases occurring at the sampling stations were recorded.

By combining the three scenarios of disease incidence given by the biological process with the three scenarios of detection function, we generated nine different scenarios of reported Incidence for each cell  $(I_i)$ , under a multinomial process given by equation (7). For each combination scenario, the response data comprised the number of reported cases per cell (Figure 1, *I*1A to *I*3C).

Preliminary manipulation of environmental layers was done using the software QGIS,<sup>82</sup> the simulations were conducted in the statistical environment R,<sup>83</sup> and Bayesian model fitting to the simulated data was carried out using the program JAGS,<sup>84</sup> interfaced with R via the package *rjags*.<sup>85</sup>

We analysed the simulated incidence data, using each of the three models described above. We used Markov Chain Monte Carlo (MCMC) algorithms (code provided in Appendix S1) to fit each of the models to the combination of environmental and incidence data. Relatively non-informative priors where chosen for all process and observation parameters and for the cells of the map relating to the latent variable. To make this a conservative test of the methodology, we employed priors wide variances. For the coefficients of the environmental covariates we chose diffuse normal priors centred at zero, corresponding to a null hypothesis of no-effect for each covariate. For the distance decay parameter  $\sigma$  of the detection function, we adopted a uniform prior with limits 0-1000.<sup>71</sup> For parameter  $\rho$  of the covariance matrix describing spatial autocorrelation in the latent covariate, we used a gamma prior (shape = 0.1, rate = 0.1). To achieve convergence, models 1 and 2 were run for  $3 \times 10^4$ , whereas model 3 was run for  $1.2 \times 10^6$  iterations.

Means of posterior distributions with corresponding credible intervals were obtained for each model coefficient  $\hat{\beta}_k$  as well as the shape parameters of the detection function  $\hat{\sigma}$ , (only relevant for models 2 and 3). For each model and each simulated scenario, we generated spatial predictions of the expected true incidence  $\hat{N}$  and the latent



**Figure 1.** Location of active sampling sites, simulation of reported malaria reported incidence  $(I_{1A}, ..., I_{3C})$  under three scenarios of insecticide resistance (IR1, IR2, IR3) and three scenarios of reporting probability as a function of distance from health centres  $(P_A, P_B, P_C)$ .

covariate of insecticide resistance  $\widehat{IR}$ . The accuracy of each parameter in the complete set  $\theta = (k, \sigma)$  was examined by calculating its relative bias from the true underlying value, as

$$RB_{\theta} = \frac{\hat{\theta} - \theta}{|\theta|} \tag{9}$$

and by plotting the simulated versus reconstructed malaria incidence (for models 2 and 3) and between the simulated and reconstructed insecticide resistance (for model 3).

# 3 Results

The full results with posterior summaries for all model parameters are reported in the supplementary material (S2). Plots showing the relationship between the simulated and reconstructed malaria incidence and between the simulated and reconstructed insecticide resistance are also presented in supplementary material (S3). Here, we present an overview of these detailed results, by reporting on the values of relative bias |RB| for each explanatory variable, in each model, under the nine different scenarios of reported malaria incidence (Figure 2).

Model 1 considered only the active sampling points, hence the single column under model 1 in Figure 2 does not include extended results pertaining to the clinic detection function (see supplementary material S2.1 for full results). Under model 1, the simulated malaria incidence was affected only by the environmental covariates (that were common to all scenarios) including insecticide resistance. Overall, the results from model 1 showed an average |RB| = 0.11 (std. dev. = 0.08). This was a persistent finding across all three simulated patterns for the latent variable (IR), with low values of relative bias arising regardless of the degree of spatial autocorrelation of the simulated insecticide resistance layer.

![](_page_7_Figure_7.jpeg)

**Figure 2.** Visual summary of results of the three Bayesian models of reported malaria incidence (*I*) under different simulated scenarios of insecticide resistance spatial patterns (*IR*) and probability of reporting at health centres (*P*). Model 1 used only active sampling data from some localised surveys, model 2 only passive case detections at health centres, model 3 combined both data sources together. The colour scale refers to the absolute values of the relative bias between the simulated coefficients of the variables involved in the biological process (1 to 6), or the shape parameter of the detection function (7), and the estimate of the same coefficient obtained by the mean of Markov Chain Monte Carlo (MCMC) posterior distributions. (•) indicates that the simulated coefficient is within the corresponding 95% posterior credible interval, (×) indicates that it falls outside.

Model 2, which considered only data from passive case detection, was less able to capture the underlying effects of predictors on the reported malaria incidence (see supplementary material S2.2 for full results). The posterior means of all coefficients showed an overall average |RB| = 0.89 (std. dev. = 1.52). A pattern of increasing bias emerged in particular when considering scenarios of increasing spatial autocorrelation in the latent variable of insecticide resistance (Figure 2). Since model 2 only included the passive detection cases, the latent variable was completely missing from the list of covariates. In scenarios *I*1A, *I*1B and *I*1C, given by the same *IR*1, (low spatial autocorrelation), the average |RB| was 0.87 (std. dev. = 1.53). Scenarios that assumed an intermediate level of spatial autocorrelation in insecticide resistance (latent variable *IR*2) generated an average |RB| of 0.89 (std. dev. = 1.54) whereas models assuming the most spatially autocorrelated distribution of insecticide resistance (*IR*3) generated an average |RB| of 0.91 (std. dev. = 1.50). Contrary to the coefficients of the process model, posteriors pertaining to the observation model were not sensitive to the different shapes of the detection function (cases  $P_A$ ,  $P_B$  or  $P_C$ ). Posteriors for the parameter  $\hat{\sigma}$  of the detection function were highly accurate, with absolute values of relative biases ranging from 0.06 to 0.09 (Figure 2). This model was able to partly reconstruct disease incidence, but not in areas with relatively higher levels of insecticide resistance (Figure 3(a)).

Model 3 gave the best results in terms of estimating coefficients with low relative biases (see supplementary material S2.3 for full results). Of particular note is the fact that the parameter for the latent insecticide resistance variable  $RB_{\hat{\beta}_{IR}}$  showed a low |RB| varying between 0.02 and 0.08. Overall, the average |RB| across all variables was 0.07 (std. dev. = 0.07). As with model 1, but in contrast to model 2, the magnitude of bias in estimated parameters was unrelated to the degree of spatial autocorrelation assumed in the latent variable. Similar to model 2, the parameter associated with the case detection function ( $\hat{\sigma}$ ) was estimated with good accuracy, but model 3 was more

![](_page_8_Figure_3.jpeg)

**Figure 3.** Reconstructions of a simulated scenario of (a) malaria incidence and (b) insecticide resistance using Bayesian models. Figure refers to scenario 3B (see Figure 1), with a high level of insecticide resistance spatial autocorrelation and an intermediate shape of the detection function. Model 1 used only active sampling data from a small set of localised surveys, model 2 only passive case detections at health centres, model 3 combined both data sources together.

accurate in mapping case distribution (Figure 3(a), see also comparison of plots in supplementary materials S3.1 vs. S3.2). Additionally, the latent distribution of the layer of insecticide resistance was accurately reconstructed using model 3 (Figure 3(b), and supplementary material S3.3).

# 4 Discussion

By analysing a wide range of plausible, simulated data sets of disease incidence and environmental variables arising from active sampling and passive case detection, we uncovered some of the disadvantages of analysing these two data types in isolation. Additionally, we propose a novel modelling framework aimed at achieving complementarity between the two. We found that such an integrated, spatially explicit model, which acknowledges both active sampling and passive case detection, leads to great improvements in precision and accuracy but also enables the reconstruction of maps for the hidden variable across unsurveyed space.

As expected, when modelling data arising only from active sampling, we achieved high explanatory power and relatively low bias, because the model had access to measurements of all the covariates underlying disease incidence. The model considering only data coming from passive case detection allowed us to estimate the map of malaria incidence with moderate accuracy. However, posterior distributions for most parameters were biased which was likely due to missing data for the important variable of insecticide resistance. This condition reflects a common situation in epidemiological studies, where passive case detection at health centres can provide a large amount of long-term data with relatively moderate effort. Our simultaneous estimation of detection functions as part of model inferences how to take account of imperfect reporting which is an integral characteristic of such opportunistic data.<sup>12,27–29,35</sup>

With our proposed third model, we achieved a good synergy between depth and breadth in inference by combining the strengths of the first two models, and allowing them to compensate for each other's limitations. In contrast to model using only passive case detection, our hybrid modelling framework allowed us to investigate the effect of all the variables (including the latent one), and to produce accurate predictive maps of the disease incidence and latent variable which were not possible with the model considering only active sampling. An important achievement of our proposed model was the capability to deal with a latent variable, regardless of its level of spatial autocorrelation. Thus, even in the absence of assumptions or any preliminary information on the spatial structure of the latent variable (e.g. whether it is akin to uncorrelated "background noise" or has a highly geography-dependent distribution), this model framework has potential to reconstruct it.

Our incremental approach showed that the gains in the accuracy of the results, moving from model 1 to model 3, were a direct result of increases in the spatial complexity used by the analytical approaches. Model 1 had no explicit spatial component. Model 2 was used to generate predictions in space but it didn't explicitly consider spatial structure in its formulation. Model 3, by including the spatial autocorrelation structure in the partly latent variable, led to the best results.

Our approach to latent variables readily generalises to processes other than insecticide resistance. We chose this particular example of a latent variable, because *IR* has potential to impact the transmission and control of a wide range of vector-borne diseases, including malaria, but is typically labour-intensive, time-consuming and expensive to measure.<sup>86</sup> Although WHO guidelines classify insecticide resistance in a binary way,<sup>86</sup> the raw data from Tube test bioassays measure the percentage survival of cohorts of similarly aged females after a given time period of exposure to insecticide resistance as a continuous variable ranging from 0 to 1. Our approach can be easily extended to more specific measures of insecticide resistance, such as metabolic, cuticular and behavioural resistance,<sup>53,87</sup> or to other types of predictor data that can be collected in the field through active sampling but are not easily obtainable via passive case detection, such as vector abundance and density.

When simulating and modelling the latent variable, we made an assumption of stationarity (the autocorrelation function did not change in space or in time) and monotonicity (the autocorrelation always decreased with distance). These two assumptions can be plausibly relaxed extending our autocorrelation function. For example, non-stationary formulations could be achieved by expressing the rate of autocorrelation decay ( $\rho$ ) as a function of latitude and longitude or time. Alternatively,  $\rho$  could be expressed as log-linear combination of environmental covariates. Non-monotonic formulations of the autocorrelation function could be produced for cases where periodic patterns exist in space, but we currently see very little justification for such formulations based on biological first principles.

The ability to account for reporting bias of our response variable makes our approach easily applicable to other scenarios where an imperfect detection needs to be considered, such as citizen science data<sup>88</sup> or mobile phone surveillance tools.<sup>89</sup> When modelling the detection function, we made similar assumptions (stationarity and

monotonicty) to those of the autocorrelation function for the latent variable and we hypothesised the observation process was only affected by distance from health centres.<sup>10–12,30,32,34</sup> In several real-world scenarios, additional covariates of reporting probability may be involved, such as age and sex of the patient and socioeconomic factors.<sup>31,33,35</sup> Borrowing fundamental concepts from Distance sampling,<sup>47</sup> we assumed that at zero distance the probability of reporting the disease was 100%; however, asymptomatic disease in apparently healthy people is common,<sup>66,90</sup> and would not be observed in clinical data. Thus, incomplete detection at zero distance (based on additional calibration data on the frequency of asymptomatic cases) must be considered.<sup>91</sup> Finally, human mobility is unlikely to be strictly related to Euclidean distance (a third implicit assumption of our detection function), so it may be preferable to use the distance according to road network,<sup>6</sup> when applying this model to real data. Global digital layers describing the travel time between any two points on the globe (based on data such as road density, terrain morphology and an political borders)<sup>36</sup> could be easily included in a spatially explicit epidemiological model such as ours.<sup>92</sup> For all of these reasons, we suggest that preliminary analysis using pilot data and focussing only on modelling the detection probability should be carried out before integrating it into the final model.

Our likelihood could be deployed using either a Bayesian or a frequentist setting. It is likely that in real life, most epidemiological data sets will be accompanied by sufficient expert opinion to lead to influential priors, hence we have illustrated using a Bayesian approach. However, we did not assume the existence of expert opinion here because we were seeking to construct a conservative test of our methods.

The models presented here (in particular our model 3, using both data types) require a high computational effort (see supplementary material for details). Notwithstanding their theoretical simplicity, the need to take spatial structure into account with a large dataset slows down the Bayesian MCMC inference. Other model fitting approaches such as the Integrated Nested Laplace approximation (INLA)<sup>93</sup> may prove capable of providing similarly accurate results but with faster processing.<sup>94</sup>

In quantitative ecology, data simulation, by generating random realisations from stochastic processes described by a series of distributional statements, is exceedingly useful.<sup>71</sup> Although simulated studies are not guaranteed to be the same as a real epidemiological system, they allow objective validation of proposed frameworks on a wide range of plausible scenarios, easily adaptable to other epidemiological studies. Although our simulation was borrowing its settings from a study specifically looking at malaria, we demonstrated its applicability on a broad range of contrasting scenarios. Therefore, we believe that such a framework can successfully work under different epidemiological systems, where a combination of large-scale but opportunistic data are collected at the same time as conducting a small number of localised scientific surveys.

The strength of our proposed analytical approach lies in its ability to use distinct solutions, such as latent process modelling and spatially autoregressive modelling, in a fully integrated framework. In particular, we demonstrated how active sampling and passive case detection, that have so far been considered independently in the context of spatial epidemiology, can be used simultaneously and complimentarily in a package where the strength of one compensates for the drawback of the other. Our method shows promise for complex spatial epidemiology studies, by allowing different parts of the model to glean information from different types of data. Such egalitarian and complementary use of two or more data types can be extended to make use of digital or hard copy primary care records, irrespective of the sophistication of the health provision systems, the density of the human population, or the nature of the disease.

## Acknowledgements

We would like to thank Laurie Baker, Julie Barker, Fergus Chadwick, Heather McDevitt and Hilary Ranson for their useful comments on an earlier draft of the manuscript. We also thank the two anonymous reviewers who provided essential suggestions to improve the presentation of the statistical elements.

## **Declaration of conflicting interests**

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

#### Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: This project has received funding from the Wellcome Trust [grant no. 200222/Z/15/Z] MiRA.

## ORCID iD

Luca Nelli D https://orcid.org/0000-0001-6091-4072

## References

- 1. Pfeiffer D, Robinson TP, Stevenson M, et al. *Spatial analysis in epidemiology*. New York, NY: Oxford University Press, 2008.
- 2. Reisen WK. Landscape epidemiology of vector-borne diseases. Ann Rev Entomol 2010; 55: 461-483.
- 3. Woolhouse M. How to make predictions about future infectious disease risks. *Philos Transact Royal Soc Lond B: Biol Sci* 2011; **366**: 2045–2054.
- 4. Hartemink N, Vanwambeke SO, Purse BV, et al. Towards a resource-based habitat approach for spatial modelling of vector-borne disease risks. *Biol Rev* 2015; **90**: 1151–1162.
- 5. Lawson AB, Banerjee S, Haining RP, et al. Handbook of spatial epidemiology. Boca Raton, FL: CRC Press, 2016.
- 6. Kirby RS, Delmelle E and Eberth JM. Advances in spatial epidemiology and geographic information systems. *Ann Epidemiol* 2017; **27**: 1–9.
- Bakuza JS, Denwood MJ, Nkwengulila G, et al. Estimating the prevalence and intensity of Schistosoma mansoni infection among rural communities in Western Tanzania: the influence of sampling strategy and statistical approach. *PLOS Neglect Tropical Dis* 2017; 11: e0005937.
- Eisen L and Eisen RJ. Using geographic information systems and decision support systems for the prediction, prevention, and control of vector-borne diseases. *Ann Rev Entomol* 2011; 56: 41–61.
- 9. Samadoulougou S, Maheu-Giroux M, Kirakoya-Samadoulougou F, et al. Multilevel and geo-statistical modeling of malaria risk in children of Burkina Faso. *Parasite Vector* 2014; 7: 350.
- 10. Müller I, Smith T, Mellor S, et al. The effect of distance from home on attendance at a small rural health centre in Papua New Guinea. *Int J Epidemiol* 1998; **27**: 878–884.
- 11. Feikin DR, Nguyen LM, Adazu K, et al. The impact of distance of residence from a peripheral health facility on pediatric health utilisation in rural western Kenya. *Tropical Med Int Health* 2009; **14**: 54–61.
- 12. Biswas RK and Kabir E. Influence of distance between residence and health facilities on non-communicable diseases: an assessment over hypertension and diabetes in Bangladesh. *PLoS ONE* 2017; **12**: e0177027.
- 13. Minuzzi-Souza TTC, Nitz N, Cuba CAC, et al. Surveillance of vector-borne pathogens under imperfect detection: lessons from Chagas disease risk (mis)measurement. *Scientific Reports* 2018; 8: 151.
- 14. Tiono AB, Kangoye DT, Rehman AM, et al. Malaria incidence in children in South-West Burkina Faso: comparison of active and passive case detection methods. *PLoS One* 2014; **9**: e86936.
- 15. Bhoomiboonchoo P, Nisalak A, Chansatiporn N, et al. Sequential dengue virus infections detected in active and passive surveillance programs in Thailand, 1994–2010. *BMC Public Health* 2015; **15**: 250.
- Sarti E, L'Azou M, Mercado M, et al. A comparative study on active and passive epidemiological surveillance for dengue in five countries of Latin America. Int J Infect Dis 2016; 44: 44–49.
- 17. Pham Thi KL, Briant L, Gavotte L, et al. Incidence of dengue and chikungunya viruses in mosquitoes and human patients in border provinces of Vietnam. *Parasite Vectors* 2017; **10**: 556.
- Fauver JR, Weger-Lucarelli J, Fakoli III LSIII, et al. Xenosurveillance reflects traditional sampling techniques for the identification of human pathogens: a comparative study in West Africa. PLOS Neglect Tropical Dis 2018; 12: e0006348.
- Mai VQ, Mai TTX, Tam NLM, et al. Prevalence and risk factors of dengue infection in Khanh Hoa Province, Viet Nam: a stratified cluster sampling survey. J Epidemiol 2018; 28: 488–497.
- Bagi J, Grisales N, Corkill R, et al. When a discriminating dose assay is not enough: measuring the intensity of insecticide resistance in malaria vectors. *Malaria J* 2015; 14: 210.
- Krajacich BJ, Slade JR, Mulligan RF, et al. Sampling host-seeking anthropophilic mosquito vectors in West Africa: comparisons of an active human-baited tent-trap against gold standard methods. *Am J Tropical Med Hygiene* 2015; 92: 415–421.
- 22. Cevallos V, Ponce P, Waggoner JJ, et al. Zika and Chikungunya virus detection in naturally infected Aedes aegypti in Ecuador. *Acta Tropica* 2018; **177**: 74–80.
- Das AK, Harries AD, Hinderaker SG, et al. Active and passive case detection strategies for the control of leishmaniasis in Bangladesh. *Public Health Action* 2014; 4: 15–21.
- 24. Insaf TZ and Talbot T. Identifying areas at risk of low birth weight using spatial epidemiology: a small area surveillance study. *Prevent Med* 2016; **88**: 108–114.
- 25. Midega JT, Smith DL, Olotu A, et al. Wind direction and proximity to larval sites determines malaria risk in Kilifi District in Kenya. *Nat Commun* 2012; **3**: 674.
- Vollack K, Sodoudi S, Névir P, et al. Influence of meteorological parameters during the preceding fall and winter on the questing activity of nymphal Ixodes ricinus ticks. *Int J Biometeorol* 2017; 61: 1787–1795.
- 27. Gething PW, Noor AM, Gikandi PW, et al. Improving imperfect data from health management information systems in Africa using space-time geostatistics. *PLOS Med* 2006; **3**: e271.

- 28. Dickersin K and Chalmers I. Recognizing, investigating and dealing with incomplete and biased reporting of clinical research: from Francis Bacon to the WHO. *J Royal Soc Med* 2011; **104**: 532–538.
- 29. Smyth RMD, Kirkham JJ, Jacoby A, et al. Frequency and reasons for outcome reporting bias in clinical trials: interviews with trialists. *BMJ* 2011; **341**: c7153.
- 30. Nemet GF and Bailey AJ. Distance and health care utilization among the rural elderly. *Social Sci Med* 2000; **50**: 1197–1208.
- 31. Kiwanuka SN, Ekirapa EK, Peterson S, et al. Access to and utilisation of health services for the poor in Uganda: a systematic review of available evidence. *Transact Royal Soc Tropical Med Hygiene* 2008; **102**: 1067–1074.
- 32. Schoeps A, Gabrysch S, Niamba L, et al. The effect of distance to health-care facilities on childhood mortality in rural Burkina Faso. *Am J Epidemiol* 2011; **173**: 492–498.
- 33. Kizito J, Kayendeke M, Nabirye C, et al. Improving access to health care for malaria in Africa: a review of literature on what attracts patients. *Malaria J* 2012; **11**: 55.
- Larson PS, Mathanga DP, Campbell CH, et al. Distance to health services influences insecticide-treated net possession and use among six to 59 month-old children in Malawi. *Malaria J* 2012; 11: 18.
- 35. Oduro AR, Maya ET, Akazili J, et al. Monitoring malaria using health facility based surveys: challenges and limitations. *BMC Public Health* 2016; **16**: 354.
- Weiss DJ, Nelson A, Gibson HS, et al. A global map of travel time to cities to assess inequalities in accessibility in 2015. *Nature* 2018; 553: 333.
- 37. Gaudart J, Touré O, Dessay N, et al. Modelling malaria incidence with environmental dependency in a locality of Sudanese savannah area, Mali. *Malaria J* 2009; **8**: 61.
- Wayant NM, Maldonado D, de Arias AR, et al. Correlation between normalized difference vegetation index and malaria in a subtropical rain forest undergoing rapid anthropogenic alteration. *Geospatial Health* 2010; 4: 179–190.
- 39. Palaniyandi M. The role of Remote Sensing and GIS for spatial prediction of vector-borne diseases transmission: a systematic review. J Vector Borne Dis 2012; 49: 197.
- 40. Parham PE, Pople D, Christiansen-Jucht C, et al. Modeling the role of environmental variables on the population dynamics of the malaria vector Anopheles gambiae sensu stricto. *Malaria J* 2012; **11**: 271.
- Hirve S, Singh SP, Kumar N, et al. Effectiveness and feasibility of active and passive case detection in the visceral leishmaniasis elimination initiative in India, Bangladesh, and Nepal. Am J Tropical Med Hygiene 2010; 83: 507–511.
- 42. Kuznetsov VN, Grjibovski AM, Mariandyshev AO, et al. A comparison between passive and active case finding in TB control in the Arkhangelsk region. *Int J Circumpolar Health* 2014; **73**: 23515.
- 43. Zhou G, Afrane YA, Malla S, et al. Active case surveillance, passive case surveillance and asymptomatic malaria parasite screening illustrate different age distribution, spatial clustering and seasonality in western Kenya. *Malaria J* 2015; 14: 41.
- Pava Z, Handayuni I, Trianty L, et al. Passively versus actively detected malaria: similar genetic diversity but different complexity of infection. Am J Tropical Med Hygiene 2017; 97: 1788–1796.
- 45. Illian J, Penttinen A, Stoyan H, et al. *Statistical analysis and modelling of spatial point patterns*. Hoboken, NJ: John Wiley & Sons, 2008.
- 46. Wiegand T and Moloney KA. Handbook of spatial point-pattern analysis in ecology. Boca Raton, FL: CRC Press, 2013.
- 47. Buckland ST. Introduction to distance sampling: estimating abundance of biological populations. Oxford: Oxford University Press, 2001.
- 48. Chandler RB and Royle JA. Spatially explicit models for inference about density in unmarked or partially marked populations. *Ann Appl Stat* 2013; 7: 936–954.
- Ramsey DS, Caley PA and Robley A. Estimating population density from presence-absence data using a spatially explicit model. J Wildlife Manage 2015; 79: 491–499.
- Bhatt S, Weiss DJ, Cameron E, et al. The effect of malaria control on *Plasmodium falciparum* in Africa between 2000 and 2015. *Nature* 2015; **526**: 207.
- 51. Gething PW, Casey DC, Weiss DJ, et al. Mapping *Plasmodium falciparum* mortality in Africa between 1990 and 2015. *New Engl J Med* 2016; **375**: 2435–2445.
- 52. Protopopoff N, Van Bortel W, Speybroeck N, et al. Ranking malaria risk factors to guide malaria control efforts in African highlands. *Plos ONE* 2009; **4**: e8022.
- Ranson H, N'Guessan R, Lines J, et al. Pyrethroid resistance in African anopheline mosquitoes: what are the implications for malaria control? *Trends Parasitol* 2011; 27: 91–98.
- Toé KH, Jones CM, N'Fale S, et al. Increased pyrethroid resistance in malaria vectors and decreased bed net effectiveness, Burkina Faso. *Emerg Infect Dis* 2014; 20: 1691–1696.
- 55. Ranson H and Lissenden N. Insecticide resistance in African Anopheles mosquitoes: a worsening situation that needs urgent action to maintain malaria control. *Trend Parasitol* 2016; **32**: 187–196.
- 56. Badolo A, Traore A, Jones CM, et al. Three years of insecticide resistance monitoring in Anopheles gambiae in Burkina Faso: resistance on the rise? *Malaria J* 2012; **11**: 232.
- 57. Govella NJ, Chaki PP and Killeen GF. Entomological surveillance of behavioural resilience and resistance in residual malaria vector populations. *Malaria J* 2013; **12**: 124.

- Guelbeogo WM, Sagnon NF, Liu F, et al. Behavioural divergence of sympatric Anopheles funestus populations in Burkina Faso. *Malaria J* 2014; 13: 65.
- 59. Mayer SV, Tesh RB and Vasilakis N. The emergence of arthropod-borne viral diseases: a global prospective on dengue, chikungunya and zika fevers. *Acta Tropica* 2017; **166**: 155–163.
- Dantas-Torres F, Chomel BB and Otranto D. Ticks and tick-borne diseases: a One Health perspective. *Trends Parasitol* 2012; 28: 437–446.
- Lai Y-S, Biedermann P, Ekpo UF, et al. Spatial distribution of schistosomiasis and treatment needs in sub-Saharan Africa: a systematic review and geostatistical analysis. *Lancet Infect Dis* 2015; 15: 927–940.
- 62. Nanyingi MO, Munyua P, Kiama SG, et al. A systematic review of Rift Valley Fever epidemiology 1931–2014. *Infect Ecol Epidemiol* 2015; **5**: 28024.
- 63. Franco JR, Simarro PP, Diarra A, et al. Epidemiology of human African trypanosomiasis. *Clin Epidemiol* 2014; 6: 257–275.
- 64. Davis JK, Vincent G, Hildreth MB, et al. Integrating environmental monitoring and mosquito surveillance to predict vector-borne disease: prospective forecasts of a West Nile virus outbreak. *PLoS Curr* 2017; **9**: 1–11.
- 65. Lawson AB. Bayesian disease mapping: hierarchical modeling in spatial epidemiology. Boca Raton, FL: CRC Press, 2013.
- 66. Bousema T, Okell L, Felger I, et al. Asymptomatic malaria infections: detectability, transmissibility and public health relevance. *Nat Rev Microbiol* 2014; **12**: 833.
- 67. Ripley BD. Spatial statistics. 1981. New York, NY: Hayward Wiley, 1981.
- Kelsall J and Wakefield J. Modeling spatial variation in disease risk: a geostatistical approach. J Am Stat Assoc 2002; 97: 692–701.
- Diboulo E, Sié A and Vounatsou P. Assessing the effects of malaria interventions on the geographical distribution of parasitaemia risk in Burkina Faso. *Malaria J* 2016; 15: 228.
- 70. Royle JA. N-mixture models for estimating population size from spatially replicated counts. Biometrics 2004; 60: 108-115.
- Kéry M and Royle JA. Applied hierarchical modeling in ecology: analysis of distribution, abundance and species richness in R and BUGS: Volume 1: Prelude and static models. Cambridge, MA: Academic Press, 2015.
- 72. Liu J and Chen X-P. Relationship of remote sensing normalized differential vegetation index to Anopheles density and malaria incidence rate. *Biomed Environ Sci: BES* 2006; **19**: 130–132.
- 73. Krefis AC, Schwarz NG, Krüger A, et al. Modeling the relationship between precipitation and malaria incidence in children from a holoendemic area in Ghana. *Am J Tropical Med Hygiene* 2011; **84**: 285–291.
- 74. Diboulo E, Sié A, Diadier DA, et al. Bayesian variable selection in modelling geographical heterogeneity in malaria transmission from sparse data: an application to Nouna Health and Demographic Surveillance System (HDSS) data, Burkina Faso. *Parasite Vectors* 2015; 8: 118.
- Srimath-Tirumula-Peddinti RCPK, Neelapu NRR and Sidagam N. Association of climatic variability, vector population and malarial disease in district of Visakhapatnam, India: a modeling and prediction analysis. *PLoS ONE* 2015; 10: e0128377.
- 76. Alout H, Roche B, Dabiré RK, et al. Consequences of insecticide resistance on malaria transmission. *PLOS Pathogen* 2017; **13**: e1006499.
- 77. Coleman M, Hemingway J, Gleave KA, et al. Developing global maps of insecticide resistance risk to improve vector control. *Malaria J* 2017; **16**: 86.
- Busetto L and Ranghetti L. MODIStsp: An R package for automatic preprocessing of MODIS Land Products time series. Comput Geosci 2016; 97: 40–48.
- 79. Silverman BW. Density estimation for statistics and data analysis. Boca Raton, FL: CRC Press, 1986.
- Gitzen RA and Millspaugh JJ. Comparison of least-squares cross-validation bandwidth options for kernel home-range estimation. *Wildlife Soc Bull* 2003; 31: 823–831.
- 81. Matowo N, Munhenga G, Tanner M, et al. Fine-scale spatial and temporal heterogeneities in insecticide resistance profiles of the malaria vector, Anopheles arabiensis in rural south-eastern Tanzania [version 1; referees: 2 approved]. 2017.
- QGIS Development Team. QGIS Geographic Information System. Open Source Geospatial Foundation Project, 2018. http://qgis.osgeo.org
- 83. R Development Core Team. R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing, 2018.
- Plummer M. JAGS: A program for analysis of Bayesian graphical models using Gibbs sampling. In: Proceedings of the 3rd international workshop on distributed statistical computing. Vienna, Austria, 2003, p.125.
- Plummer MS, Alexey Denwood, Matt. *rjags: Bayesian graphical models using MCMC*. Version 4.6, https://cran.r-project. org/web/packages/rjags/index.html (2016).
- World Health Organization. *Malaria entomology and vector control*. Geneva, Switzerland: World Health Organization, 2013.
- Corbel V and N'Guessan R. Distribution, mechanisms, impact and management of insecticide resistance in malaria vectors: a pragmatic review. In: Manguin S (Ed.) *Anopheles mosquitoes New insights into malaria vectors*. 2013, pp. 579–633.

- Robinson OJ, Ruiz-Gutierrez V and Fink D. Correcting for bias in distribution modelling for rare species using citizen science data. *Diversity Distribution* 2017; 24: 460–472.
- 89. Mtema Z, Changalucha J, Cleaveland S, et al. Mobile phones as surveillance tools: implementing and evaluating a largescale intersectoral surveillance system for rabies in Tanzania. *PLOS Med* 2016; **13**: e1002002.
- 90. Lindblade KA, Steinhardt L, Samuels A, et al. The silent threat: asymptomatic parasitemia and malaria transmission. *Expert Rev Anti-infective Ther* 2013; **11**: 623–639.
- Laake JL and Borchers DL. Methods for incomplete detection at distance zero. Adv Distance Sampl Oxford, UK: Oxford University Press, 2004, pp.108–189.
- 92. Alegana VA, Atkinson PM, Lourenço C, et al. Advances in mapping malaria for elimination: fine resolution modelling of *Plasmodium falciparum* incidence. *Scientific Rep* 2016; **6**: 29628.
- Rue H, Martino S and Chopin N. Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. J Royal Stat Soc: Ser B (Stat Methodol) 2009; 71: 319–392.
- 94. Blangiardo M and Cameletti M. Spatial and spatio-temporal Bayesian models with R-INLA. UK: John Wiley & Sons, 2015.