

1 **Modelling the impact of MAUP on environmental drivers for *Schistosoma***  
2 ***japonicum* prevalence**

3 **Andrea L. Araujo Navas <sup>1,\*</sup>, Frank Osei <sup>1</sup>, Ricardo J. Soares Magalhães <sup>2,3</sup>, Lydia R.**  
4 **Leonardo <sup>4</sup>, and Alfred Stein <sup>1</sup>**

5

6 <sup>1</sup> Faculty of Geo-information Science and Earth Observation (ITC), University of Twente,  
7 PO Box 217, 7500 AE, Enschede, The Netherlands.

8 <sup>2</sup> UQ Spatial Epidemiology Laboratory, School of Veterinary Science, The University of  
9 Queensland, Gatton 4343 QLD, Australia.

10 <sup>3</sup> Child Health and Environment Program, Child Health Research Centre, The University  
11 of Queensland, South Brisbane 4101 QLD, Australia. Institute of Biology, College of  
12 Science, University of the Philippines Diliman, the Philippines.

13

14 \* Correspondence: [a.l.araujonavas@gmail.com](mailto:a.l.araujonavas@gmail.com)

15

16 E-mails:

17 ALAN: [a.l.araujonavas@utwente.nl](mailto:a.l.araujonavas@utwente.nl)

18 FO: [f.b.osei@utwente.nl](mailto:f.b.osei@utwente.nl)

19 RJSM: [r.magalhaes@uq.edu.au](mailto:r.magalhaes@uq.edu.au)

20 LRL: [lydialeonardo1152@gmail.com](mailto:lydialeonardo1152@gmail.com)

21 AS: [a.stein@utwente.nl](mailto:a.stein@utwente.nl)

22

23 **Abstract**

24

25 **Background:** The modifiable areal unit problem (MAUP) arises when the support size of a  
26 spatial variable affects the relationship between prevalence and environmental risk factors. Its  
27 effect on schistosomiasis modelling studies could lead to unreliable parameter estimates. The  
28 present research aims to quantify MAUP effects on environmental drivers of *Schistosoma*  
29 *japonicum* infection by (i) bringing all covariates to the same spatial support, (ii) estimating  
30 individual-level regression parameters at 30 m, 90 m, 250 m, 500 m, and 1 km spatial  
31 supports, and (iii) quantifying the differences between parameter estimates using five models.

32 **Methods:** We modelled the prevalence of *Schistosoma japonicum* using sub-provinces health  
33 outcome data and pixel-level environmental data. We estimated and compared regression  
34 coefficients from convolution models using Bayesian statistics.

35 **Results:** Increasing the spatial support to 500 m gradually increased the parameter estimates  
36 and their associated uncertainties. Abrupt changes in the parameter estimates occur at 1 km  
37 spatial support, resulting in loss of significance of almost all the covariates. No significant  
38 differences were found between the predicted values and their uncertainties from the five  
39 models. We provide suggestions to define an appropriate spatial data structure for modelling  
40 that gives more reliable parameter estimates and a clear relationship between risk factors and  
41 the disease.

42 **Conclusions:** Inclusion of quantified MAUP effects was important in this study on  
43 schistosomiasis. This will support helminth control programs by providing reliable parameter  
44 estimates at the same spatial support, and suggesting the use of an adequate spatial data  
45 structure, to generate reliable maps that could guide efficient mass drug administration  
46 campaigns.

47 **Keywords:** schistosomiasis modelling; modifiable areal unit problem; uncertainty; Bayesian  
48 statistics; convolution model.

## 49 **Background**

50

51 Schistosomiasis (SCH) is a water-borne neglected tropical disease of public health  
52 significance [1] associated with important morbidity outcomes in school-aged children such  
53 as malnutrition, anaemia and stunted growth in school-aged children [2, 3]. Infection is  
54 caused by skin penetration of the cercariae, the larval infective stage of the parasite, also  
55 known as schistosome. Three schistosome species cause the infection: *Schistosoma*  
56 *japonicum*, *S.mansoni*, and *S. haematobium*. Due to its zoonotic life cycle [4], *Schistosoma*  
57 *japonicum* is the hardest to control; its infection life cycle includes the amphibious snail from  
58 the species *Oncomelania hupensis* as the intermediate host, and humans and other  
59 mammals as definite hosts [5, 6]. SCH affects more than 252 million people worldwide [7]  
60 especially populations living in poor conditions, where access to clean water and sanitation is  
61 limited.

62 Traditionally, SCH is controlled by the use of anthelmintic drugs in at-risk  
63 populations [8]. Mass drug administration campaigns identify at-risk populations by using

64 SCH risk mapping. SCH mapping uses geographic information systems (GIS), global  
65 positioning systems and remotely sensed environmental data [9, 10]. Modelling those  
66 infections using various statistical methods have enabled the study of the distribution of  
67 populations at-risk [9, 10], and the role of the environmental variation on the geographical  
68 heterogeneity of infection burden (i.e. prevalence or intensity of infection) [11]. Statistical  
69 modelling of SCH quantifies empirical relationships between indirect morbidity indicators of  
70 public health significance and environmental risk factors. Those could be extracted from  
71 Earth Observation (EO) data such as monitor sites or satellite imagery. In addition, EO data  
72 help to interpolate the level of infection towards unsampled locations [12-14].

73         The robustness of SCH geographical modelling efforts is affected by uncertainties  
74 propagated from the use of EO data at various spatial and temporal scales of analysis [15]. EO  
75 data are generally constrained by their spatial and temporal scale of sampling [16]. In this  
76 study, we focus on spatial scale. Scale is a major concern in spatial epidemiology [17, 18]  
77 since it determines the significance of the various environmental risk factors on the disease  
78 distribution [19]. Spatial scale encompasses the spatial support and the spatial extent of  
79 analysis [20]. The spatial support refers to the area that each individual observation occupies  
80 in space. In the case of a raster grid, the spatial support is the spatial resolution (e.g. a 30 x 30  
81 m-resolution Landsat pixel). The spatial extent is the spatial coverage of a set of observations  
82 (e.g. administrative units) and is gathered following a sampling scheme [20]. For a given  
83 extent, the support size and shape of spatial units may affect the patterns identified in the  
84 survey and environmental data [21, 22] and the relationship between the disease morbidity  
85 indicators and the environmental risk factors. This is known as the modifiable areal unit  
86 problem (MAUP) [23, 24]. The MAUP arises because spatial units of analysis are often  
87 created using different *ad hoc* shapes and sizes. Statistical analyses of data performed  
88 according to these varying spatial units may lead to different results (e.g. correlation and  
89 regression coefficients) [24].

90         Various studies investigated the consequences of ignoring MAUP effects in spatial  
91 epidemiological modelling. For instance, Hellsten et al. [25] studied the influence of using  
92 aggregated covariate data to model ammonia emissions at the farm level. They showed that  
93 the size and shape of spatial aggregation areas strongly affect the location of the emissions  
94 estimated by the model, e.g. too small areas resulting in false emission “hot spots”. Schur et al  
95 [21] and Schur et al [22] aggregated SCH prevalence maps to estimate endemicity for various  
96 administrative units [26]. Such aggregation showed different patterns of endemicity and

97 intervention approaches. As a consequence, localized areas of high endemicity may not be  
 98 addressed properly. In a recent study [27] we quantified the effect of pure specification bias,  
 99 that originates when using group-level (i.e. aggregated) survey data at an administrative level  
 100 for individual-level inferences. Equation 1 shows the common method used to model  
 101 schistosomiasis. Data on the human *Schistosoma japonicum* infection variable  $y$  are  
 102 commonly aggregated at barangay  $k$  level,  $y_k$  has a binomial distribution with parameters  $N_k$   
 103 and  $p_k$  corresponding to the number of sampled individuals and the probability of infection,  
 104 respectively. Parameters for this distribution are obtained from the mean of various  
 105 environmental risk factors within barangay  $k$  as predictors, denoted as  $\bar{x}_k$ , where  $\gamma$  are the  
 106 barangay-level coefficients,  $\gamma_0$  being the intercept, and  $\gamma_{(1...n)}$  the regression coefficients for  $n$   
 107 environmental covariates (equation 1).

$$y_k | \bar{x}_k, \boldsymbol{\gamma} \sim \text{Binomial}(N_k, \hat{p}_k)$$

$$\text{logit}(\hat{p}_k) = \gamma_0 + \gamma_1 \cdot \bar{x}_{1_k} + \gamma_2 \cdot \bar{x}_{2_k} + \dots + \gamma_n \cdot \bar{x}_{n_k}. \quad (1)$$

108

109 We calculated individual-level regression parameters by modifying equation 1 into a  
 110 convolution model. We observed differences ranging from -0.19 to 0.28 between individual  
 111 (i.e.  $\gamma$  coefficients) and group level parameter estimates and their uncertainties. High  
 112 differences were observed for NDWI (0.28), LSTN (-0.19) and LSTD (0.16). Although some  
 113 covariates got a less significant effect on schistosomiasis, uncertainties in their individual  
 114 level coefficients were lower than the group-level regression coefficients (e.g. LSTD and  
 115 elevation). We concluded that the choice of spatial support affects the model parameter  
 116 estimates and their associated uncertainties by changing the within covariates variability in  
 117 exposure areas. The selection of spatial support should be further investigated as it might  
 118 represent a significant source of uncertainty in SCH modelling [15].

119 Up to date, the majority of SCH studies have put little attention to the size of spatial  
 120 support. They use EO data at various spatial supports with misaligned grids ignoring the  
 121 possible consequences on the observed patterns of the data [21, 22]. Moreover, MAUP effects  
 122 on the various environmental risk factors used as drivers for SCH infection have not been  
 123 quantified. This is important as the relevance of the environmental risk factors on SCH  
 124 depends on their scale of analysis [7, 19]. Ignoring MAUP effects might produce unreliable  
 125 predictions of at-risk populations, and consequently, wrong decisions based upon inefficient  
 126 mass drug administration campaigns.

127 The purpose of this research is to quantify MAUP effects on environmental drivers of  
128 *Schistosoma japonicum* infection. To achieve this objective we aim to: (i) aggregate and  
129 disaggregate EO data in order to bring all covariates to a the same spatial support of analysis,  
130 (ii) estimate individual-level covariate regression parameters at 30 m, 90 m, 250 m, 500 m,  
131 and 1 km spatial supports, by using a convolutional model that accounts for pure specification  
132 bias; and (iii) quantify the differences between parameter estimates using five different  
133 models.

## 134 **Methods**

135

### 136 **Study Area and Data on *Human Schistosoma japonicum* infection**

137 We use *Schistosoma japonicum* infection data collected as part of the 2008 Nationwide  
138 schistosomiasis Survey in the Philippines. Here, *Schistosoma japonicum* is endemic in 28 of  
139 its 81 provinces [28], with approximately 1.8 million estimated infected people [29]. The  
140 disease affects children, adolescents, and individuals with high-risk occupations, such as  
141 farmers and fishermen [29, 30]. The area of study is the region of Mindanao in the Philippines  
142 (Figure 1). This area was selected due to the high response rate of 70.9 percent of the  
143 individuals to the 2008 survey [31, 32], and the good spatial coverage of the sampling.

144 A two-stage systematic cluster sampling was used where stratification was done using  
145 high, medium and low prevalence levels, obtained from the 1994 World Bank-assisted  
146 Philippine Health Development Program. Provinces and sub-municipalities called barangays  
147 were the primary and secondary sampling units respectively. A barangay is the smallest  
148 administrative division in the Philippines, numbering from 58 to 1158 within a single  
149 province. In total 11 provinces with high ( $\geq 2\%$ ) and medium (0.091 -1.99%) prevalence  
150 rates were included, while 9 low-prevalence (0.04-0.09%) provinces were randomly selected.  
151 Within the selected provinces, barangays with high prevalence rates were surveyed. In total,  
152 between 2 and 10 barangays were surveyed per province, resulting in 108 out of 10021  
153 barangays that were surveyed in Mindanao.

154 For *Schistosoma japonicum* diagnosis, a Kato-Katz thick smear examination [32] was  
155 used based on a two-sample stool collection. Due to inconsistencies in the second stool  
156 sample submission, however, only the results of the first sample were available [8]. Samples  
157 were taken from people aged two years and above and were analysed using a microscope.  
158 Active infection was indicated by the presence of *Schistosoma japonicum* eggs.

159 Data such as age and gender were recorded for 19763 individuals. Barangay and  
160 province information for each individual was recorded but not geo-referenced. For this  
161 reason, individual-level survey data were aggregated and geo-located to the centroids of the  
162 108 barangays. We used a probability of infection  $p$  in barangay  $k$  as our disease outcome  
163 variable. We obtained an up-to-date barangay centroids shape file from DIVA geographic  
164 information system [33]. More details about the sampling design and surveyed information  
165 can be found in Leonardo et al [31, 34].

### 166 **Environmental risk factors**

167 We included in our analysis six relevant environmental risk factors for SCH transmission [35,  
168 36]. These are the nearest distance to water bodies (NDWB), the normalized difference  
169 vegetation index (NDVI), the normalized difference water index (NDWI), land surface  
170 temperature at day (LSTD) and at night (LSTN) and elevation (E). NDWB shows the  
171 accessibility of people to water bodies that represent potential infection foci as they may  
172 contain contaminated snail hosts that release the infective larval stages of the parasite [8].  
173 NDVI is an indicator of flooded vegetation [8], particularly rice paddy fields, and  
174 environmental moisture [37, 38]. Both are an important risk factor for Asian SCH [39].  
175 NDWI was used as a proxy indicator of flooding [37, 40] showing potentially hidden water  
176 bodies. LSTD and LSTN are determinant for the survival of larval stages of snails [41, 42]  
177 and are used as proxies for water temperature given that the thermal condition of shallow  
178 waters usually reflects the ambient temperature of the air [8]. Elevation is relevant for SCH  
179 transmission as the local topography of the area determines the presence of snails [43, 44].  
180 For instance, at lower altitudes the risk of finding snails increases.

181 NDWB values range from 0.17 to 26.2 km and were calculated using the closest  
182 facility network analysis tool from ArcGIS [45]. We used the river and road network, and the  
183 cities and hamlets locations as input for the network. Rivers and roads were extracted from  
184 the Open Street Map Project in the Philippines [46]. Cities and hamlet locations were  
185 obtained from the National Mapping and Resource Information Authority from The  
186 Philippines [47] data base from 2010. We calculated the nearest distance from each city and  
187 hamlet to a water body following a road and interpolated those values within all surveyed  
188 barangays towards a spatial support of 30 m.

189 NDVI values range from 0 to 0.84 and were obtained from two sources of  
190 information, i.e. a series of Landsat 5 images from 2008 with a spatial support of 30 m and

191 the MODIS MOD13Q1 product with a spatial support of 250 m. NDWI values range from  
192 0.06 to 0.61 and were also obtained from two sources of information, i.e. a Landsat 5 imagery  
193 product from 2008 with a spatial support of 30 m and the annual composite from Landsat 7  
194 from 2008 derived from Google Earth Engine with a spatial support of 500 m. LSTD values  
195 range from 297.77 to 309.52 °K and LSTN ranges from 289.73 to 297.29 °K. LSTD and  
196 LSTN values were derived from MODIS MOD11A2\_LST product with a spatial support of 1  
197 km. Finally, Elevation values range from 0 to 969.57 m was obtained from ASTER GDEM  
198 version 2 from USGS [48] with a spatial support of 30 m. All covariates were set to a  
199 common coordinate system UTM zone 51N and were standardized to have mean = 0 and  
200 standard deviation = 1 before being used. Table 1 summarizes all sources of information.

201

### 202 *Modifying the areal units of analysis*

203 From now onwards we will refer to an areal unit as the spatial support of analysis (SSA). We  
204 used five SSAs, with a spatial support equal to 30 m, 90 m, 250 m, 500 m, and 1 km,  
205 respectively. These spatial supports increase when going from low to high data aggregation.  
206 These values were selected based upon the commonly used spatial supports at which the  
207 environmental information is originally provided.

208

209 For NDVI, SSA = 30 m, we obtained NDVI values from Landsat 5 images. Many of  
210 these images presented gaps due to the presence of clouds. These gaps were covered using  
211 disaggregated NDVI MODIS images at the Landsat resolution. Disaggregation was  
212 performed using a linear model that predicted NDVI Landsat values based on NDVI MODIS  
213 values. NDVI values were obtained by merging the original and predicted Landsat NDVI  
214 values. For SSA = 90 m, we aggregated the previously merged NDVI values using their  
215 mean. For SSA = 250 m, we used the NDVI MODIS product directly. Finally, for SSA = 0.5  
216 and 1 km, we aggregated the NDVI mean values from MODIS.

217 NDWI values were obtained from the Landsat 5 images. Gaps in some of these images  
218 were covered using disaggregated NDWI composite images at the Landsat resolution.  
219 Disaggregation towards SSA = 30 m was done by interpolating NDWI values using ordinary  
220 Kriging. For SSA = 90 m and 250 m, we aggregated the combined 30 m NDWI using its  
221 mean. For SSA = 500m, we directly used the Landsat 7 composite. Finally, for SSA = 1km,  
222 we aggregated the mean of the original Landsat 7 composite.

223 To obtain LSTD and LSTN values for SSA = 30 m we disaggregated the original  
224 MODIS values by using Ordinary Kriging interpolation. For SSA = 90 m, 250 m, and 500 m,  
225 we aggregated the previously interpolated values using their mean. For SSA =1 km, we used  
226 directly LSTD and LSTN from MODIS.

227 The interpolated NDWB values for SSA = 30 m were used to obtain NDWB for SSA  
228 = 90 m, 250 m, 500 m, and 1 km by aggregating the mean values. For elevation, we directly  
229 used the original 30 m SSA Aster images. For SSA = 90 m, 250 m, 500 m, and 1 km, we  
230 aggregated the mean values of the original Aster images.

## 231 **Modelling *Schistosoma japonicum* infection under the MAUP**

232

### *Convolution model*

233 We modelled human *Schistosoma japonicum* infection at the five increasing SSAs  
234 using a convolution model that accounts for pure specification bias [27]. Pure specification is  
235 a source of uncertainty [11, 49] that produces loss of information on the real relationship  
236 between the disease and the environmental covariate data, when using aggregated survey data  
237 in a non-linear model, for example for individual-level inferences [50]. It is called ‘pure’  
238 because it specifically addresses model specification bias [51], and it biases the estimates  
239 because any direct link between exposure and health outcomes is imperfectly measured [52].  
240 This is because the regression function does not approximate the real relationship between the  
241 affected population and their exposure [27]. Pure specification bias can be reduced as the  
242 within area exposure is more homogenous [50]. This could be done by having a finer partition  
243 of space at which environmental risk factors are available [50, 53].

244 In this study we propose to minimize and quantify pure specification bias by  
245 extracting covariate information from cities within barangays (Figure 2) and by modelling the  
246 disease using a convolution model [53]. The city level is the finest available extent of  
247 analysis. Cities thus serve as a proxy for individual-level exposure locations. We identified all  
248 cities within the surveyed barangays using Google Earth Images. Available cities were  
249 extracted from the 2010 build-up data base from the National Mapping and Resource  
250 Information Authority from the Philippines [54]. We completed unavailable cities using  
251 Google Earth Images.



252 For the convolution model, we used the aggregate data method proposed by Prentice  
 253 and Sheppard [55]. For each SSA, we obtained covariate information  $\mathbf{x}$  for image pixel  $i$   
 254 belonging to a city  $j$  within a specific barangay  $k$  (Figure 2). Let  $n = 6$  be the number of  
 255 covariates  $\mathbf{x}_{ijk}$  measured at locations  $s_{ijk}$   $i = 1, \dots, m_k$  where  $m_k$  denotes the number of city  
 256 pixels within barangay  $k$ . Note that with an increasing resolution the possibility increases that  
 257 there are no pixel points falling in cities of the within-pixel sizes. Data on the human  
 258 *Schistosoma japonicum* infection are available at individual-level recorded within a barangay  
 259  $k$ . Because the exact response locations of the individual-level data are unknown, we  
 260 aggregated them to their corresponding barangay centroid, denoted by  $y_k$ . To estimate the  
 261 average probability of infection of the individuals in barangay  $k$ , and the individual level  
 262 coefficients  $\boldsymbol{\beta}$ , we obtained the mean risk function  $\hat{p}_k$  over the total number of city pixels or  
 263 exposure locations (equation 2). We accounted for the spatial variability at barangay-level by  
 264 adding a spatial structure random effects term  $s_k$ . Pure specification bias results as  $\gamma \neq \beta$ , and  
 265 is then minimized by using the individual-level regression coefficients  $\boldsymbol{\beta}$  instead of the group-  
 266 level coefficients  $\boldsymbol{\gamma}$ . The accompanying uncertainties are quantified by the difference between  
 267 the group and individual-level credible intervals [27] for each SSA. The convolution model  
 268 used is of the form:

$$y_k | \mathbf{x}_{ijk}, \boldsymbol{\beta} \sim \text{Binomial}(N_k, \hat{p}_k)$$

$$\hat{p}_k = \frac{1}{m_k} \cdot \sum_{i=1}^{m_k} \frac{1}{1 + \exp(-(\beta_0 + \beta_1 \cdot x_{1ijk} + \beta_2 \cdot x_{2ijk} + \dots + \beta_n \cdot x_{nijk} + s_k))}. \quad (2)$$

## 269 Model implementation

270 Five models were implemented, all including an intercept ( $\beta_0$ ), pixel-level environmental  
 271 variables ( $\mathbf{x}_{ijk}$  = NDVI, NDWI, LSTD, LSTN, E, NDWB) and their corresponding individual-  
 272 level coefficients  $\boldsymbol{\beta}$ . Collinearity between covariates was assessed with the Pearson  
 273 correlation coefficient. All covariates were standardized to have mean = 0 and standard  
 274 deviation = 1.

275 The intercept  $\beta_0$  was given a diffuse uniform prior distribution with wide bounds  $\beta_0$   
 276  $\sim U[-100, 100]$ . The other  $\boldsymbol{\beta}$  parameters were given a diffuse normal distribution  
 277  $\boldsymbol{\beta} \sim N[0, \frac{1}{\sigma^2}]$ , with  $\sigma$  uniformly distributed on a wide range of  $\sigma \sim U[0, 100]$ . These  
 278 distributions avoid overestimating the parameters [56] and allow a good sequences mixing

279 used for Markov Chain Monte Carlo (MCMC) simulations, contributing to a fast convergence  
280 [57].

281 Prior information for the spatially structured random effects was based upon a geo-  
282 statistical model that can be used as a sampling distribution for continuous spatial data [58].  
283 The vector of random variables  $\mathbf{s}$  associated with point locations  $(x_k, y_k)$ ,  $k = 1, \dots, K$ , was  
284 modelled with a multivariate normal distribution  $s \sim MVN_K[\mu, \Sigma_{ab}]$  with mean  $\mu = 0$  and a  
285 covariance matrix  $\Sigma_{ab} = \sigma^2 \cdot \exp[-(\phi \cdot d_{ab})^\kappa]$  defined by a powered exponential spatial  
286 decaying correlation function.

287 The covariance matrix  $\Sigma_{ab}$  is specified as a function of the distances  $d_{ab}$  between  
288 barangay centroids  $a$  and  $b$ , with the rate of decline of spatial correlation per unit of distance  
289  $\phi$ , the scalar parameter representing the overall variance  $\sigma^2$  and the scalar parameter  $\kappa$   
290 controlling the amount of spatial smoothing. Because extreme values of  $\kappa$  (0 and 2) could lead  
291 to undesirable smoothing, we used  $\kappa = 1$ . Prior information for  $\phi$  was set to be  
292 uniform:  $\phi \sim U[2E10^{-7}, 3E10^{-3}]$ . These values give a diffuse but plausible prior range of  
293 correlations between 0.1 and 0.99 at the minimum distance between points (575 m), and  
294 between 0 and 0.3 at the maximum distance between points ( $< 552$  km), assisting  
295 identifiability [59]. For  $\sigma^2$ , a half-normal distribution was selected:  $\sigma^2 \sim HN[0,1]$  to restrict  
296 the prior  $\sigma^2$  to positive values and avoid problems with convergence [56, 60].

297 To run the model we used three sequences or chains with 50000 iterations. This  
298 number of iterations ensured that the simulations were representative of target distributions  
299 and a stable convergence [57]. In order to diminish the influence of starting values, we  
300 discarded the first half of each sequence [57] using a burn-in of 25000 iterations.  
301 Convergence was monitored visually and statistically by inspecting the trace plots, and by  
302 checking the  $\hat{R}$  statistic [61, 62] also called the potential scale reduction factor. This potential  
303 scale reduction factor assesses sequences mixing by comparing the between and within  
304 variation. An  $\hat{R}$  value  $< 1.1$  indicates evidence that sequences have converged [61], while  
305 higher values suggest that an increase in the number of simulations may improve the  
306 inferences [57].

307 Survey and environmental data were structured in a rectangular format where columns  
308 are headed by the array name. Survey data and the codes in BUGS for the various SSA are  
309 provided in the Additional files 1 and 2 respectively.

310

### 311 **Model validation**

312 The five models were validated using two methods. The first method compared the data  
313 generated from the simulations of the predictive distribution to the observed data using a test  
314 statistic. A test statistic is a value derived from the sampled data and is used to perform  
315 hypothesis testing. This test statistic is the posterior predictive  $p$ -value (ppp-value) generated  
316 by calculating the proportion of the predicted values which are more extreme than the  
317 observed maximum, minimum and mean prevalence observed value. We calculated (i) the  
318 proportion of simulations of the data from the model for which the maximum prevalence  
319 across simulated barangays is greater than or equal to the maximum observed value, (ii) the  
320 proportion of simulations of the data from the model for which the minimum prevalence  
321 across simulated barangays is greater than or equal to the minimum observed value, and (iii)  
322 the proportion of simulations of the data from the model for which the mean prevalence  
323 across simulated barangays is greater than or equal to the mean observed value. If the model  
324 fits the data, the simulated values distribute closely around the observed values, thus, we  
325 expect a ppp-value of around 0.5. Otherwise, for a biased model, the ppp-value will be close  
326 to 0 or 1. Our aim was to test whether the model predicts a similar number of barangays with  
327 maximum and minimum prevalence values compared with the observed data. We generated  
328 ppp values for maximum, minimum and mean prevalence values for the models at five  
329 increasing SSA using 75000 simulations. The second method used the area under the curve  
330 (AUC) of the receiving operating characteristics (ROC). We applied a threshold of 0.5%  
331 (prevalence mean in Mindanao region) since we are interested in knowing the ability of the  
332 models to discriminate the mean prevalence level in the study area. We also examined the  
333 ability of the model to discriminate the number of positive cases, thus, we used a threshold of  
334 1, which indicates the presence of at least one positive case. We used an AUC value of 70%  
335 to indicate acceptable predictive performance [8, 63].

### 336 **Software**

337 Model implementation was done in the software OpenBUGS 3.2.3 [64, 65] (Medical research  
338 Council, Cambridge, UK, and Imperial College London, UK). It was downloaded for free at  
339 [66]. We called Open BUGS from R using the package R2OpenBUGS [67]. The spatial  
340 models were coded using the GeoBUGS [59] function as an add-on module to OpenBUGS.  
341 GeoBUGS provides an interface to work with conditional autoregressive and geo-statistical  
342 models. Data pre-processing and Ordinary Kriging was performed in R [68].

## 343 **Results**

### 344 **Modelling *Schistosoma japonicum* infection under the MAUP**

345

#### 346 *Convolution model*

347 Our findings show that NDVI has a non-significant effect on the prevalence of SCH infection  
348 for all SSA, except for SSA = 1 km (Table 2 and Figure 3a). NDVI estimates vary gradually  
349 from 0.19 to 0.26 when increasing SSA until 500 m. For SSA = 1 km, the estimate rapidly  
350 increases to 0.59. Uncertainties are similar throughout all SSA (Figure 4a and Table 3),  
351 slightly increasing when increasing SSA. The highest credible interval value is 0.60 for SSA  
352 = 250 m, and the lowest is 0.52 for SSA = 30 m.

353 NDWI has a significant negative effect on the prevalence of SCH infection throughout  
354 all SSAs (Table 2 and Figure 3b). When SSA increases, parameter estimates increase from -  
355 1.06 to -0.76, coming somewhat closer to zero. We found similar estimates for SSA = 30 m,  
356 90 m and 250 m (i.e. -1.06 to -1.02), and for SSA = 500 m and 1 km (i.e. -0.8 to -0.76)  
357 (Figure 4b). Uncertainty values are similar for all SSAs and show a slight decrease when  
358 increasing SSA (Figure 4b and Table 3). The highest uncertainty value equals 0.54 for SSA =  
359 30 m and the lowest value equals 0.44 for SSA = 500 m.

360 LSTD has a significant negative effect on the prevalence of SCH infection for almost  
361 all SSA, except for SSA = 1 km (Table 2 and Figure 3c). Similar parameter estimates equal to  
362 -0.71 are obtained for SSA = 30 m, 90 m, and 250 m, while the parameter estimate increases  
363 slightly to -0.65 for SSA = 500 m. For SSA = 1 km, there is a noticeable increase in the  
364 parameter estimate to -0.01 (Figure 3c and Figure 4c). Uncertainty increases from 0.59 to 0.64  
365 when increasing SSA from 30 m to 500 m, but for SSA = 1 km there is a considerable  
366 increase in uncertainty to 1.49 (Figure 4c).

367 LSTN has a significant negative effect on the prevalence of SCH infection for almost  
368 all SSA, except for SSA = 1 km (Table 2 and Figure 3d). Parameter estimates increase from -  
369 0.78 to -0.86 while increasing SSA from 30 m to 500 m. For SSA = 1 km, the parameter  
370 estimate rapidly goes up to 0.1 (Figure 3d and Figure 4d). Uncertainty increases slightly from  
371 0.56 to 0.58 when increasing SSA from 30 m to 500 m, but it increases considerably to 1.14  
372 for SSA = 1 km (Table 3 and Figure 4d).

373 Elevation has a significant negative effect on the prevalence of SCH infection for all  
374 SSA, except for SSA = 1 km (Table 2 and Figure 3e). When increasing SSA from 30 m to

375 500 m, parameter estimates slightly decrease from -0.95 to -1.03. For SSA = 1 km, the  
376 parameter estimate considerably increases to -0.04 (Figure 3e and 4e). Uncertainty values  
377 vary from 0.59 to 0.64 when increasing SSA from 30 m to 500 m. For SSA = 1 km,  
378 uncertainty considerably decreases to 0.35 (Table 3 and Figure 4). The lowest uncertainty  
379 value is 0.35 for SSA = 1 km and the highest is 0.66 for SSA = 250 m.

380 Finally, NDWB has a significant negative effect on the prevalence of SCH infection  
381 for all SSA (Table 2 and Figure 3f). We found similar parameter estimates of -0.28, -0.29 and  
382 -0.31 for SSA = 30 m, 90 m, and 250 m, respectively, and estimates of -0.38 and -0.4 for SSA  
383 = 500 m and 1 km, respectively (Figure 3f). Uncertainties constantly increase from 0.32 to  
384 0.39 (Table 3 and Figure 4f) when increasing SSA.

385 Intercept values range from -6.02 to -6.17 for almost all SSAs, except for SSA = 1 km,  
386 where it is equal to -5.49. The rate of decay of spatial autocorrelation ( $\phi$ ) ranges from  
387  $1.65 \times 10^{-5}$  to  $2.81 \times 10^{-4}$  for SSAs = 1 km and 500 m, respectively.

388 Our findings show high and moderate correlation and determination ( $R^2$ ) coefficient values  
389 between the SSAs and all environmental covariates estimates (Table 4) with correlation  
390 coefficients ranging from -0.94 to 0.94 and  $R^2$  values from 0.6 to 0.86, respectively.

391 Correlation coefficients between the SSAs and uncertainties are high for LSTD, LSTN and  
392 NDWB with the values of 0.91, 0.9 and 0.91, respectively. Determination coefficients  $R^2$   
393 between the SSAs and uncertainties in the covariates estimates are moderate for LSTD, LSTN  
394 and NDWB with the values of 0.76, 0.75 and 0.76, respectively (Table 4). Uncertainties in  
395 NDVI and NDWI estimates do not show any correlation with SSAs (Table 4).

396 **Table 2.** *Regression coefficient estimates for each risk factor at five descending spatial*  
397 *supports of analysis*

### 398 ***Influence on predictions***

399 Differences between observed and predicted prevalence values are similar for the five SSA  
400 models (Figure 5). Variation in these differences is highest between the 30 m and 1 km  
401 models ( $R^2 = 0.94$ ) and lowest between the 30 m and 90 m models ( $R^2 = 0.99$ ). Figure 5  
402 shows that the maximum and minimum differences are 1.11% and 0.01%, respectively,  
403 corresponding to the 1 km SSA model. For fitted prevalence values higher than 2% all models  
404 underestimate the prevalence of infection, while, for fitted prevalence values lower than 2%,  
405 overestimation and underestimation occur for the five models (Figure 5). A plot of the

406 residuals against prevalence from Figure 5 serve as a visual inspection of the fit, where we  
407 realize that it is based on positive predictive predictions.

408           Uncertainties on the predictions are similar for the five models (Additional file 3).  
409 Higher differences in uncertainty were found between the 500 m and 1 km models ( $R^2 =$   
410  $0.96$ ), and lower differences were found between the 90 m and 250 m models ( $R^2 = 0.99$ ).  
411 The highest uncertainty value is 9.23% for all the models, except the 1 km model with 8.9%,  
412 and the lowest uncertainty value is 0.006% for the 1 km model.

### 413 **Model validation**

414 The maximum and minimum observed prevalence values are 8.5% and 0.33%, respectively.  
415 The first validation method shows *ppp*-values for all SSA ranging from 0.64 to 0.67 for the  
416 first test statistic (Table 5). This means that simulated data slightly deviate from around 0.14  
417 to 0.17 from the maximum observed prevalence value (Figure 6). For all SSA it is likely to  
418 see a similar number of predicted maximum prevalence values compared to the observed data.  
419 For the second test statistic, *ppp*-values ranged from 0.87 to 0.93 (Table 5). This means that  
420 simulated data are biased around 0.36 to 0.43 from the minimum observed prevalence data  
421 (Figure 7). For almost all SSA, simulated data predict a higher number of minimum  
422 prevalence values compared to the observed data. For the last test statistics, *ppp*-values  
423 ranged from 0.59 to 0.67 (Table 5), showing that simulated data deviate from around 0.09 to  
424 0.17 from the mean observed prevalence value (Figure 8).

425           Results from the second validation method show that all models have a high ability to  
426 predict prevalence values, with AUC values of 0.91 for SSA = 30 m, 90 m, 250 m, and 500  
427 m, and 0.93 for SSA = 1 km. All models have a good ability to predict the positive number of  
428 SCH cases. Models with SSA = 30 m, 90 m, 250 m, and 500 m models have AUC values of  
429 0.83, while the 1 km SSA model presents a lower AUC value of 0.79, showing a decrease in  
430 the ability to predict the positive number of SCH cases.

### 431 **Discussion**

432 Schistosomiasis modelling studies have commonly used environmental risk factors as drivers  
433 for disease exposure and transmission [69, 70]. The studies so far have used spatially  
434 misaligned environmental variables at different spatial supports of analysis, ignoring MAUP  
435 effects on the parameter estimates, predictions, and the relationship between disease  
436 morbidity indicators and risk factors. This study is the first effort to quantify the effects of  
437 modifying the areal unit (i.e. spatial support) of NDVI, NDWI, LSTD, LSTN, E, and NDWB,

438 on model parameter estimates and their uncertainties. Uncertainty may be quantified using  
439 measures of accuracy or imprecision [15]. We evaluated uncertainty using measures of  
440 imprecision based on the nature of the disease and the survey data available and quantified it  
441 using credible intervals in a Bayesian setting. We applied it to *Schistosoma japonicum*  
442 infection modelling in the Mindanao region, the Philippines.

443 Our findings show that the environmental risk factors NDVI, NDWI, LSTD, LSTN,  
444 and E behave similarly when increasing the SSA from 30 m to 1 km (Table 4). An increase in  
445 SSA from 30 m to 500 m does not represent any significant changes in parameter estimates.  
446 Conversely, for SSA = 1km, all show a considerable increase in their estimates. The reasons  
447 are explained below.

448 NDVI has a positive effect on SCH, meaning that higher NDVI values increase the  
449 prevalence of infection. This is explained by the positive relationship between vegetation,  
450 moisture and snail density [37]. NDVI effects are not significant for SSA < 1 km, because  
451 NDVI is an indicator of greenness that is mainly effective for arid areas and Mindanao is not  
452 arid. However, the NDVI effect becomes significant on the prevalence of SCH infection for  
453 SSA = 1 km. This could be because NDVI effects on SCH prevalence are greater at global  
454 scales [8] than at local scales. This might be explained by the fact that prevalence values at  
455 local scales can vary significantly at nearby locations, as it depends not only on the nature of  
456 the parasite lifecycle, which requires optimal habitat conditions (i.e. environmental  
457 conditions), but also on sanitation conditions on the area [71]. The increase in uncertainty  
458 values with increasing SSA is due to the coarse areal pixels  $\geq 250$  m resolution that does not  
459 reliably represent rice paddy fields. Those are substantially smaller than 25 ha, i.e. are  
460 covered by at most four pixels [72].

461 For SSA = 30 m, 90 m, 250 m, and 500 m, LSTD, LSTN, and E have a significant  
462 negative effect on SCH prevalence. Conversely for SSA = 1 km, their parameter estimates are  
463 close to zero. This means that when the areal unit reaches 1 km, the effect of these covariates  
464 on the prevalence of SCH infection becomes non-significant. This is also observed from the  
465 credible intervals of these covariates for the 1 km SSA model. The reason is that the  
466 homogeneity of the covariate values increases when increasing the SSA. LST is a proxy of the  
467 ambient temperature of the air, which reflects the thermal conditions of shallow waters [27].  
468 Its negative relationship with the prevalence of infection could be explained by the fact that  
469 temperatures above 19 to 20°C do not influence the release of cercariae from the infected host

470 to the infection foci [73], as well as temperatures below approximately 15°C arrest the  
471 development of *S.japonicum* in the snail host [74]. The minimum LST value at night is  
472 around 21°C, while the maximum LST value during the day is 31°C. LSTD and LSTN  
473 uncertainty values for SSA = 1 km are remarkably high as compared to other SSA. This is  
474 explained by the coarse LSTD and LSTN areal pixels of 1 km<sup>2</sup> that cannot reliably represent  
475 low and high temperature zones in city areas that range from 0.02 to 3 km<sup>2</sup> [27]. Elevation has  
476 a negative effect on SCH. This was expected as in areas with high elevation values (>2300m)  
477 the risk of infection is low [75]. Conversely, the risk of infection is high for elevation areas  
478 below 900m. Elevation uncertainty values are similar for all SSA, except for SSA = 1 km,  
479 where its value considerably decreases to 0.34. Here we see the effect of the gradual changes  
480 of elevation in Mindanao region are gradual and without steep slopes [27]. Using data directly  
481 at the 1 km SSA could give reliable elevation values, but with a non-significant effect on the  
482 disease prevalence.

483 For NDWB and NDWI, an increase in SSA from 30 m to 250 m represents non-  
484 significant changes in parameter estimates, which range from -1.06 to -1.02 for NDWI, and  
485 from -0.31 to -0.28 for NDWB. Conversely, when increasing the SSA to 500 m, parameter  
486 estimates change to -0.8 and -0.38 for NDWI and NDWB, respectively. For SSA = 500 m and  
487 1 km, NDWI estimates increase, having a less significant effect on SCH prevalence, again  
488 due to the increase in the homogeneity of the covariate values when increasing SSA. Higher  
489 NDWI values show the presence of potential hidden infection foci. Nevertheless, results show  
490 that NDWI presents a negative effect on SCH (Table 2). This could be because NDWI cannot  
491 efficiently suppress the signal from build-up land mixing enhanced water features with build-  
492 up land noise. Thus, build-up noise could also have high NDWI values [40]. According to Gu  
493 et al. [76] NDWI values lower than 0.3 indicate the presence of drought areas. In our study  
494 area, we found that around 77% of Mindanao present drought conditions, explaining the  
495 negative effect on the disease. NDWB estimates decrease when increasing SSAs (Table 4),  
496 specially for SSA = 500 m and 1 km, but their significance on SCH prevalence increases. A  
497 possible explanation is that people that move larger distances to water bodies are most likely  
498 to get infected. This could be because at spatial supports < 1 km, NDWB values seem to be  
499 more homogenous than at smaller spatial supports, showing a weaker relationship with the  
500 disease (Table 2). For spatial supports > 1km, neighbouring pixels present more  
501 heterogeneous values, which could be because of the aggregation process, but also because of  
502 the use of some kind of transportation media that allows apparent reduction of travel distances



503 in a relatively large area (1 km<sup>2</sup>). Clearly, transportation (type of road and media of  
504 transportation) plays an important role [77].

505           Uncertainty values for NDWI decrease when increasing the SSA, with a minimum of  
506 0.44 for SSA = 500 m. Clearly, NDWI data originally available at SSA = 250 m are more  
507 reliable than values modified to larger SSAs. Using Ordinary Kriging for interpolation  
508 increases the variance in the estimates in a somewhat unrealistic way since it uses a constant  
509 mean [58], while in reality, means are different. Uncertainty values of NDWB, for instance,  
510 increase with increasing SSA due to the coarse areal pixel units  $\geq 0.25$  km<sup>2</sup>. Such a size is  
511 insufficient to reliably define nearest distances to water bodies in city areas of 0.02 to 3 km<sup>2</sup>.

512           Our aim was not to compare the performance of the models as we used the same  
513 model structure, number and type of covariates in the five models. Thus the model itself is  
514 practically the same for all SSA. Although our aim was not focused on model comparison, the  
515 resulting DIC values from Table 2 suggest the use of spatial support sizes below or equal to  
516 250 m in SCH modelling. This is shown by the low DIC values from 86.67 to 140.5 for SSA  
517  $\leq 250$  m, and high DIC values of 143.7 and 147.5 for the 500 m and 1 km models,  
518 respectively (Table 2).

519           When modelling prevalence of *Schistosoma japonicum* infection in Mindanao, the  
520 effect of increasing SSA, or modifying the areal unit of analysis, from 30 m to 500 m,  
521 produces a gradual and continuous increase on the parameter estimates and their associated  
522 uncertainties. For SSA = 1 km, sudden changes occur in the relationship between the risk  
523 factors and the prevalence of the disease. This is shown by the non-significant effect of almost  
524 all explanatory variables on *Schistosoma japonicum* prevalence. Results suggest that the use  
525 of environmental data extracted at SSA = 1 km is not appropriate for the modelling of  
526 *Schistosoma japonicum* prevalence.

527           Bayesian statistical methods were used to model the disease, and along with a  
528 convolution regression model, they corrected for pure specification bias on our estimates.  
529 This is a relevant contribution to the analysis of uncertainties in this type of spatial  
530 epidemiological study. For future studies, new trends in geospatial artificial intelligence  
531 (geoAI), that could resolve limitations regarding the MAUP for exposure modelling studies,  
532 are emerging to model schistosomiasis [78] as well as other diseases [79]. We particularly  
533 identified (i) the use of high-performance computing to handle spatiotemporal big data, and  
534 (ii) machine and deep learning algorithms implementation to big data infrastructures to extract

535 relevant disease or environmental information [79, 80]. One example is a data-driven method  
536 used to predict particulate matter air pollution ( $PM_{2.5}$ ) in Los Angeles, CA, USA. Here,  
537 machine learning was used on spatial big data, i.e. land use and roads, derived from  
538 OpenStreetMap, to predict  $PM_{2.5}$  concentrations. When generating relative importance  
539 measures for the different risk factors, MAUP effects reduced when applying a random forest  
540 model that was trained with the distances between the features and the monitoring  $PM_{2.5}$   
541 stations, [81]. The rapid development of geoAI methods, their advantage to deal with big data,  
542 and their rapid computational time, makes them an attractive and advantageous tool to tackle  
543 limitations with modelling schistosomiasis and other diseases. There is still little work done in  
544 this field, but we think it is valuable to further explore geoAI solutions to deal with the  
545 MAUP, and perhaps other inherent uncertainties produced in disease modelling and mapping.

546         Finding MAUP effects on the various environmental risk factors used for modelling  
547 *Schistosoma japonicum* prevalence, is a step forward to the uncertainty analysis in the  
548 schistosomiasis, and possibly other diseases. The present research deals with limitations such  
549 as the use of aggregated disease data, due to the lack of geo-located individual-level surveys.  
550 It also provides a robust method for the selection of an appropriate spatial data structure,  
551 which at the same time, enables the acquisition of more reliable parameter estimates, and  
552 defines a clear relationship between the risk factors and the disease. From the public health  
553 perspective, this research can support helminth control programs by providing less uncertain  
554 models and maps. Epidemiologists and health scientists could use these maps to identify risk  
555 areas for the control and prevention of the disease [12, 82], which in the case of  
556 schistosomiasis, is generally based on mass drug administration campaigns addressed to the  
557 identified at-risk populations. The provision of reliable information is relevant to guide mass  
558 drug administration campaigns by enhancing the assessment of the infection risk,  
559 understanding its potential impacts on human health [15, 83] and avoiding erroneous  
560 conclusions and decisions about the spatial distribution of schistosomiasis [15, 27]. This  
561 research is also relevant to evaluate the effectiveness of mass drug administration campaigns,  
562 as it could guide the identification of persistent hotspots, or places where prevalence of  
563 infection remains despite mass drug administration efforts [71]. It is known that despite the  
564 implementation of mass drug administration campaigns, some places do not show a decrease  
565 in local SCH transmission. This is because these campaigns do not only depend on the nature  
566 of the parasite lifecycle and the poor sanitation conditions, but also on the local environmental

567 factors, drivers for SCH transmission. Finding relevant environmental factors at local level  
568 would allow more intensive efforts at persistent hotspots.

## 569 **Conclusions**

570 The present study shows a clear MAUP effect on *Schistosoma japonicum* modelling. An  
571 increase in parameter estimates and their associated uncertainties occurs when increasing the  
572 spatial support of analysis (SSA). It also showed that using environmental data extracted at  
573 SSA = 1 km is not relevant for *Schistosoma japonicum* prevalence of infection at this specific  
574 extent of analysis, as this leads to wrong conclusions about the distribution of the disease and  
575 its relationship with the potential risk factors. Thus, the use of maps based upon this  
576 information is to be avoided as these may guide health scientists in the control or prevention  
577 of the disease astray. The results from this study could guide other disease modelling studies  
578 as they suggest a spatial support sizes at which environmental information has no longer a  
579 significant effect on the disease, and which data structure is recommended for the modelling.  
580 Epidemiologists, decision makers, and health scientists could thus benefit from those, e.g. to  
581 better understand and quantify MAUP effects on the relationship between the disease and its  
582 risk factors, and to provide reliable maps that are useful for disease control and prevention.

## 583 **Additional files**

584 **Additional file 1: Table S1.** Survey data aggregated at the barangay level. These data show  
585 the number of positive cases ( $y$ ) and the total number of sampled people ( $n$ ) in a barangay ( $k$ ).

586 **Additional file 2: Text S1. Code** for the convolution model used in OpenBUGS. It includes  
587 the prior distributions used for the covariate and spatial parameters, and the model itself. The  
588 model uses two indexes,  $k$  for the barangay and  $j$  for the number of city pixels within  
589 barangays. The betas are the individual-level regression coefficients,  $m$  is the number of city  
590 pixels in a barangay. Covariates are  $ndvi$ : normalized difference vegetation index,  $ndwi$ :  
591 normalized difference water index,  $lstd$ : land surface temperature day,  $lstn$ : land surface  
592 temperature night,  $e$ : elevation,  $ndwb$ : nearest distance to water bodies. The spatial parameter  
593 is represented as  $s$ .

594 **Additional file 3: Figure S1.** Residual plot for the five increasing spatial supports of analysis.  
595 a) SSA = 30m, b) SSA = 90m, c) SSA= 250m, d) SSA= 500m, e) SSA= 1km. The x axis  
596 represents the fitted prevalence values for the five spatial supports of analysis. The y axis

597 represents the residuals calculated by the difference between the observed and predicted  
598 prevalence values.

## 599 **Abbreviations**

600 ASTER: Advanced spaceborne thermal emission and reflection radiometer; AUC: Area under  
601 the curve; BUGS: Bayesian inference using Gibbs sampling; E: Elevation; EO: Earth  
602 observation; geoAI: Geospatial artificial intelligence; GDEM: global digital elevation map;  
603 GIS: Geographic information systems; LSTD: Land surface temperature day; LSTN: Land  
604 surface temperature night; MAUP: Modifiable areal unit problem; MCMC: Markov chain  
605 Monte Carlo; MODIS: Moderate resolution imaging spectroradiometer; NDVI: Normalized  
606 difference vegetation index; NDWB: Nearest distance to water bodies; NDWI: Normalized  
607 difference water index; ROC: Receiver operator characteristic; SCH: Schistosomiasis; SSA:  
608 Spatial support of analysis; USGS: United states geological survey.

609

## 610 **Declarations**

### 611 **Ethics Approval**

612 The data used in this study were collected in 2005 when there was no requirement for ethical  
613 review and clearance. This study used aggregated survey data at the barangay level, which  
614 enabled the full de-identification of individuals involved in the survey.

615 The study results represent part of the PhD thesis entitled “Statistical evaluation of spatial  
616 uncertainty in schistosomiasis mapping”, which was published as: Araujo Navas  
617 AL. Statistical evaluation of spatial uncertainty in schistosomiasis mapping. Enschede:  
618 University of Twente, Faculty of Geo-Information Science and Earth Observation (ITC),  
619 2019. 157 p. (ITC Dissertation). <https://doi.org/10.3990/1.9789036548281>.

### 620 **Consent for publication**

621 Not applicable

### 622 **Availability of data and material**

623 The datasets used and/or analysed during the current study are available from the  
624 corresponding author on reasonable request.

### 625 **Competing interests**

626 The authors declare that they have no competing interests

## 627 **Funding**

628 The research team acknowledges the World Health Organization funds provided to conduct  
629 the surveys. This research received no external funding. ALAN's doctoral research is funded  
630 by the University of Twente. The funders had no role in study design, data collection, and  
631 analysis, decision to publish, or preparation of the manuscript.

## 632 **Authors' contributions**

633 Conceptualization, ALAN, FO and RJSM; Formal analysis, ALAN; Investigation, ALAN;  
634 Methodology, ALAN and FO; Resources, LRL and RJSM; Software, ALAN; Supervision,  
635 FO, RJSM and AS; Validation, ALAN; Visualization, ALAN; Writing – original draft,  
636 ALAN; Writing – review & editing, ALAN, FO, RJSM and AS. All authors read and  
637 approved the final

## 638 **Acknowledgments**

639 Not applicable

## 640 **References**

641

- 642 1. Walz Y, Wegmann M, Dech S, Vounatsou P, Poda J-N, N'Goran EK, et al. Modeling and  
643 validation of environmental suitability for schistosomiasis transmission using remote sensing. *PLoS*  
644 *Negl Trop Dis*. 2015;9:e0004217.
- 645 2. Leenstra T, Acosta LP, Langdon GC, Manalo DL, Su L, Olveda RM, et al. *Schistosoma*  
646 *japonicum*, anemia, and iron status in children, adolescents, and young adults in Leyte, Philippines.  
647 *Am J Clin Nutr*. 2006;83:371–79.
- 648 3. Coutinho HM, McGarvey ST, Acosta LP, Manalo DL, Langdon GC, Leenstra T, et al. Nutritional  
649 status and serum cytokine profiles in children, adolescents, and young adults with *Schistosoma*  
650 *japonicum*-associated hepatic fibrosis, in Leyte, Philippines. *J Infect Dis*. 2005;192:528–36.
- 651 4. Jia TW, Zhou XN, Wang XH, Utzinger J, Steinmann P, Wu XH. Assessment of the age-specific  
652 disability weight of chronic *Schistosoma japonicum*. *Bull World Health Organ*. 2007;85:458–65.
- 653 5. Tarafder MR, Balolong E, Carabin H, Belisle P, Tallo V, Joseph L, et al. A cross-sectional study  
654 of the prevalence of intensity of infection with *Schistosoma japonicum* in 50 irrigated and rain-fed  
655 villages in Samar province, the Philippines. *Bmc Public Health*. 2006;6:10.
- 656 6. Yang K, Wang XH, Yang GJ, Wu XH, Qi YL, Li HJ, et al. An integrated approach to identify  
657 distribution of *Oncomelania hupensis*, the intermediate host of *Schistosoma japonicum*, in a  
658 mountainous region in China. *Int J Parasitol*. 2008;38:1007–16.
- 659 7. Hotez PJ, Alvarado M, Basanez MG, Bolliger I, Bourne R, Boussinesq M, et al. The global  
660 burden of disease study 2010: interpretation and implications for the neglected tropical diseases.  
661 *PLoS Negl Trop Dis*. 2014;8:e2865.
- 662 8. Soares Magalhães RJ, Salamat MS, Leonardo L, Gray DJ, Carabin H, Halton K, et al.  
663 Geographical distribution of human *Schistosoma japonicum* infection in The Philippines: tools to  
664 support disease control and further elimination. *Int J Parasitol*. 2014;44:977–84.
- 665 9. Herbreteau V, Salem G, Souris M, Hugot J-P, Gonzalez J-P. Thirty years of use and  
666 improvement of remote sensing, applied to epidemiology: from early promises to lasting frustration.  
667 *Health Place*. 2007;13:400–3.

- 668 10. Kalluri S, Gilruth P, Rogers D, Szczur M. Surveillance of arthropod vector-borne infectious  
669 diseases using remote sensing techniques: a review. *PLoS Pathogens*. 2007;3:e116.
- 670 11. Zhang ZJ, Manjourides J, Cohen T, Hu Y, Jiang QW. Spatial measurement errors in the field of  
671 spatial epidemiology. *Int J Health Geogr*. 2016;15:12.
- 672 12. Soares Magalhães RJ, Clements ACA, Patil AP, Gething PW, Brooker S. The applications of  
673 model-based geostatistics in helminth epidemiology and control. *Adv Parasitol*. 2011;74:267–96.
- 674 13. Cadavid Restrepo AM, Yang YR, McManus DP, Gray DJ, Giraudoux P, Barnes TS, et al. The  
675 landscape epidemiology of echinococcoses. *Infec Dis Poverty*. 2016;5.
- 676 14. Weiss DJ, Mappin B, Dalrymple U, Bhatt S, Cameron E, Hay SI, et al. Re-examining  
677 environmental correlates of *Plasmodium falciparum* malaria endemicity: a data-intensive variable  
678 selection approach. *Malar J*. 2015;14.
- 679 15. Araujo Navas AL, Hamm NAS, Soares Magalhães RJ, Stein A. Mapping soil transmitted  
680 helminths and schistosomiasis under uncertainty: a systematic review and critical appraisal of  
681 evidence. *PLoS Negl Trop Dis*. 2016;10:e0005208.
- 682 16. Wang XH, Zhou XN, Vounatsou P, Chen Z, Utzinger J, Yang K, et al. Bayesian spatio-temporal  
683 modeling of *Schistosoma japonicum* prevalence data in the absence of a diagnostic 'gold' standard.  
684 *PLoS Negl Trop Dis*. 2008;2:e250.
- 685 17. Walz Y, Wegmann M, Leutner B, Dech S, Vounatsou P, N'Goran EK, et al. Use of an  
686 ecologically relevant modelling approach to improve remote sensing-based schistosomiasis risk  
687 profiling. *Geospat Health*. 2015;10:271–79.
- 688 18. Young LJ, Gotway CA, Yang J, Kearney G, DuClos C. Assessing the association between  
689 environmental impacts and health outcomes: a case study from Florida. *Stat Med*. 2008;27:3998–  
690 4015.
- 691 19. Simoonga C, Utzinger J, Brooker S, Vounatsou P, Appleton CC, Stensgaard AS, et al. Remote  
692 sensing, geographical information system and spatial analysis for schistosomiasis epidemiology and  
693 ecology in Africa. *Parasitology*. 2009;136:1683–93.
- 694 20. Atkinson PM, Graham AJ. Issues of scale and uncertainty in the global remote sensing of  
695 disease. In: Hay SI, Graham A, Rogers DJ, editors. *Global mapping of infectious diseases: methods,  
696 examples and emerging applications*. San Diego: Elsevier Academic Press Inc; 2006. p. 79–118
- 697 21. Schur N, Hurlimann E, Garba A, Traore MS, Ndir O, Ratard RC, et al. Geostatistical model-  
698 based estimates of schistosomiasis prevalence among individuals aged ≤ 20 Years in west Africa.  
699 *PLoS Negl Trop Dis*. 2011;5:e1194.
- 700 22. Schur N, Hurlimann E, Stensgaard AS, Chimfwembe K, Mushingi G, Simoonga C, et al.  
701 Spatially explicit schistosoma infection risk in eastern Africa using bayesian geostatistical modelling.  
702 *Acta Trop*. 2013;128:365–77.
- 703 23. Dungan JL, Perry JN, Dale MRT, Legendre P, Citron-Pousty S, Fortin MJ, et al. A balanced view  
704 of scale in spatial statistical analysis. *Ecography*. 2002;25:626–40.
- 705 24. Openshaw S. *The modifiable areal unit problem*. : Norwick: GeoBooks; 1984.
- 706 25. Hellsten AS. *A spatio-temporal ammonia emissions inventory for the UK*.  
707 <https://www.era.lib.ed.ac.uk/handle/1842/24693>. Accessed April 2018.
- 708 26. Schur N, Vounatsou P, Utzinger J. Determining treatment needs at different spatial scales  
709 using geostatistical model-based risk estimates of schistosomiasis. *PLoS Negl Trop Dis*. 2012;6:e1773.
- 710 27. Araujo Navas AL, Osei F, Leonardo LR, Soares Magalhães RJ, Stein A. Modeling *Schistosoma*  
711 *japonicum* Infection under pure specification bias: impact of environmental drivers of infection. *Int J*  
712 *Env Res Pub He*. 2019;16:176.
- 713 28. Leonardo L, Rivera P, Saniel O, Solon JA, Chigusa Y, Villacorte E, et al. New endemic foci of  
714 schistosomiasis infections in the Philippines. *Acta Trop*. 2015;141:354–60.
- 715 29. Leonardo L, Acosta LP, Olveda RM, Aligui GDL. Difficulties and strategies in the control of  
716 schistosomiasis in the Philippines. *Acta Trop*. 2002;82:295–99.
- 717 30. Zhou XN, Bergquist R, Leonardo L, Yang GJ, Yang K, Sudomo M, et al. *Schistosoma japonicum*:  
718 control and research needs. In: Zhou XN, Bergquist R, Olveda R, Utzinger J, editors. *Important*

- 719 helminth infections in southeast Asia: diversity and potential for control and elimination, Part A. San  
720 Diego Elsevier Academic Press Inc; 2010. p. 145–78.
- 721 31. Leonardo L, Rivera P, Sanieel O, Villacorte E, Leaban MA, Crisostomo B, et al. A national  
722 baseline prevalence survey of schistosomiasis in the Philippines using stratified two-step systematic  
723 cluster sampling design. *J Trop Med*. 2012;2012:8.
- 724 32. Leonardo LR, Rivera P, Sanieel O, Villacorte E, Crisostomo B, Hernandez L, et al. Prevalence  
725 survey of schistosomiasis in Mindanao and the Visayas, The Philippines. *Parasitol Int*. 2008;57:246–  
726 51.
- 727 33. DIVA-GIS free, simple and effective. Hijmans R, Rojas E, Cruz M, O'Brien R, Barrantes I.,  
728 University of California and International Potato Center in Peru 2018. <http://www.diva-gis.org/Data>.  
729 Accessed 8 April 2018.
- 730 34. Santos FLN, Cerqueira E, Soares NM. Comparison of the thick smear and Kato-Katz  
731 techniques for diagnosis of intestinal helminth infections. *Rev Soc Bras Med Trop*. 2005;38:196–98.
- 732 35. Brooker S, Hay S, Issae W, Hall A, Kihamia C, Lawambo N, et al. Predicting the distribution of  
733 urinary schistosomiasis in Tanzania using satellite sensor data. *Trop Med Int Health*. 2001;6:998–  
734 1007.
- 735 36. Brooker S, Hay SI, Tchuente L-AT, Ratard R. Using NOAA-AVHRR data to model human  
736 helminth distributions in planning disease control in Cameroon, West Africa. *Photogramm Eng Rem  
737 S*. 2002;68:175–79.
- 738 37. Walz Y, Wegmann M, Dech S, Raso G, Utzinger J. Risk profiling of schistosomiasis using  
739 remote sensing: approaches, challenges and outlook. *Parasit Vectors*. 2015;8:16.
- 740 38. Malone JB, Yilma JM, McCarroll JC, Erko B, Mukaratirwa S, Zhou XY. Satellite climatology and  
741 the environmental risk of *Schistosoma mansoni* in Ethiopia and east Africa. *Acta Trop*. 2001;79:59–  
742 72.
- 743 39. Zhou YB, Liang S, Jiang QW. Factors impacting on progress towards elimination of  
744 transmission of *Schistosoma japonicum* in China. *Parasit Vectors*. 2012;5:7.
- 745 40. Xu H. Modification of normalised difference water index (NDWI) to enhance open water  
746 features in remotely sensed imagery. *Int J Remote Sens*. 2006;27:3025–33.
- 747 41. Woolhouse M, Chandiwana S. Population dynamics model for *Bulinus globosus*, intermediate  
748 host for *Schistosoma haematobium*, in river habitats. *Acta Trop*. 1990;47:151–60.
- 749 42. Pietrock M, Marcogliese DJ. Free-living endohelminth stages: at the mercy of environmental  
750 conditions. *Trends Parasitol*. 2003;19:293–99.
- 751 43. Stensgaard AS, Jorgensen A, Kabatereine NB, Rahbek C, Kristensen TK. Modeling freshwater  
752 snail habitat suitability and areas of potential snail-borne disease transmission in Uganda. *Geospat  
753 Health*. 2006;1:93-104.
- 754 44. Stensgaard AS, Utzinger J, Vounatsou P, Hurlimann E, Schur N, Saarnak CFL, et al. Large-scale  
755 determinants of intestinal schistosomiasis and intermediate host snail distribution across Africa: does  
756 climate matter? *Acta Trop*. 2013;128:378–90.
- 757 45. ESRI. ArcGIS desktop. New release simplifies your work, provides new ways to share  
758 information, supplies GIS in the cloud. 10 ed: Environmental systems research institute; 2011.
- 759 46. Planet OSM. Open Street Map Project. 2017. <https://planet.osm.org>. Accessed 21 Nov 2017.
- 760 47. Humanitarian data exchange v1.25.3. United nations office for the coordination of  
761 humanitarian affairs, New York and Geneva. 2018.  
762 [https://data.humdata.org/search?groups=phi&q=&ext\\_page\\_size=25](https://data.humdata.org/search?groups=phi&q=&ext_page_size=25). Accessed 10 April 2018.
- 763 48. Global Data Explorer. U.S. Department of interior 2017. <https://gdex.cr.usgs.gov/gdex/>.  
764 Accessed 7 August 2017.
- 765 49. King G. A solution to the ecological inference problem: reconstructing individual behavior  
766 from aggregate data. 1st ed. Princeton: Princeton University Press; 2013.
- 767 50. Wakefield J, Lyons H. Spatial aggregation and the ecological fallacy. In: G. F, editor. Handbook  
768 of modern statistical methods. Boca Raton: Chapman & Hall/CRC Press; 2010. p. 541–58.

- 769 51. Gelfand AE, Diggle P, Guttorp P, Fuentes M. Handbook of spatial statistics. 1st ed. Boca  
770 Raton: Taylor & Francis Group; 2010.
- 771 52. Richardson S, Monfort C. Ecological correlation studies. In: Elliot P, Wakefield JC, Best NG,  
772 Briggs DJ, editors. Spatial epidemiology: methods and applications. Oxford: Oxford University Press;  
773 2000. p. 205–20.
- 774 53. Wakefield J, Shaddick G. Health-exposure modeling and the ecological fallacy. *Biostatistics*.  
775 2006;7:438–55.
- 776 54. National Mapping and Resource Information Authority (NAMRIA). Department of  
777 Environment and Natural Resources. 2018. <http://www.namria.gov.ph/>. Accessed 3 February 2018.
- 778 55. Prentice RL, Sheppard L. Aggregate data studies of disease risk factors. *Biometrika*.  
779 1995;82:113–25.
- 780 56. Gelman A. Prior distributions for variance parameters in hierarchical models (comment on an  
781 article by Browne and Draper). *Bayesian Anal*. 2006;1:515–33.
- 782 57. Gelman A, Carlin JB, Stern HS, Dunson D, Rubin DB. Bayesian data analysis. 3rd ed. New York:  
783 Taylor & Francis Group; 2013.
- 784 58. Diggle PJ, Tawn J, Moyeed R. Model-based geostatistics. *J R Stat Soc Ser C Appl Stat*.  
785 2002;47:299–350.
- 786 59. Thomas A, Best N, Lunn D, Arnold R, Spiegelhalter D. GeoBugs user manual. In: MRC  
787 Biostatistics Unit. 2004. <https://www.mrc-bsu.cam.ac.uk/software/bugs/thebugs-project-geobugs/>.  
788 Accessed February 2018.
- 789 60. Lunn D, Jackson C, Best N, Thomas A, Spiegelhalter D. The BUGS book: a practical  
790 introduction to bayesian analysis. 1st ed. Boca raton: Taylor & Francis Group; 2012.
- 791 61. Brooks SP, Gelman A. General methods for monitoring convergence of iterative simulations. *J*  
792 *Comput Graph Stat*. 1998;7:434–55.
- 793 62. Gelman A, Rubin DB. Inference from iterative simulation using multiple sequences. *Stat Sci*.  
794 1992;7:457–72.
- 795 63. Brooker S, Hay SI, Bundy DAP. Tools from ecology: useful for evaluating infection risk  
796 models? *Trends Parasitol*. 2002;18:70–4.
- 797 64. Spiegelhalter D, Thomas A, Best N, Lunn D. WinBUGS user manual. In: MRC Biostatistics Unit.  
798 2003. <http://www.mrc-bsu.cam.ac.uk/bugs>. Accessed February 2018.
- 799 65. Spiegelhalter D, Thomas A, Best N, Lunn D. OpenBUGS user manual, version 3.0. 2. In: MRC  
800 Biostatistics Unit. 2007. <http://www.openbugs.net/w/Manuals>. Accessed February 2018.
- 801 66. Lunn D. SD, Thomas A. and Best N. OpenBUGS version 3.0.2. In: Downloads. 2018.  
802 <http://www.openbugs.net/w/Downloads>. Accessed 25 June 2018.
- 803 67. Sturtz S, Ligges U, Gelman A. R2OpenBUGS: a package for running OpenBUGS from R. *Journal*  
804 *of Statistical Software*. 2005;12:1–16.
- 805 68. Team RDC. R: A language and environment for statistical computing. Vienna, Austria: The R  
806 foundation for statistical computing; 2013.
- 807 69. Hu Y, Bergquist R, Lynn H, Gao F, Wang Q, Zhang S, et al. Sandwich mapping of  
808 schistosomiasis risk in Anhui province, China. *Geospat Health*. 2015;10:111–16.
- 809 70. Stensgaard AS, Vounatsou P, Sengupta ME, Utzinger J. Schistosomes, snails and climate  
810 change: current trends and future expectations. *Acta Trop*. 2019;190:257–68.
- 811 71. Kittur N, Binder S, Campbell Jr CH, King CH, Kinung'hi S, Olsen A, et al. Defining persistent  
812 hotspots: areas that fail to decrease meaningfully in prevalence after multiple years of mass drug  
813 administration with praziquantel for control of schistosomiasis. *Am J Trop Med Hyg*. 2017;97:1810–  
814 17.
- 815 72. Rice science for a better world. In: International Rice Research Institute. 2018.  
816 <http://irri.org/our-work/research/policy-and-markets/mapping-rice-in-the-philippines-where>.  
817 Accessed 22 October 2018.
- 818 73. Bauman PM, Bennett HJ, Ingalls Jr JW. The Molluscan intermediate host and *Schistosoma*  
819 *japonicum*. *Am J Trop Med Hyg*. 1948;1:567–75.



- 820 74. Yang G-J, Utzinger J, Sun L-P, Hong Q-B, Vounatsou P, Tanner M, et al. Effect of temperature  
821 on the development of *Schistosoma japonicum* within *Oncomelania hupensis*, and hibernation of *O.*  
822 *hupensis*. *Parasitology Research*. 2007;100:695–700.
- 823 75. Pesigan TP, Hairston NG, Jauregui JJ, Garcia EG, Santos AT, Santos BC, et al. Studies on  
824 *Schistosoma japonicum* infection in The Philippines 2. The molluscan host. *Bull World Health Organ*.  
825 1958;18:481–578.
- 826 76. Gu Y, Brown JF, Verdin JP, Wardlow B. A five-year analysis of MODIS NDVI and NDWI for  
827 grassland drought assessment over the central great plains of the United States. *Geophys Res Lett*.  
828 2007;34.
- 829 77. Kumm M, de Moel H, Ward PJ, Varis O. How close do we live to water? A global analysis of  
830 population distance to freshwater bodies. *PLoS One*. 2011;6:e20578.
- 831 78. Mari L, Gatto M, Ciddio M, Dia ED, Sokolow SH, De Leo GA, et al. Big-data-driven modeling  
832 unveils country-wide drivers of endemic schistosomiasis. *Sci Rep*. 2017;7:11.
- 833 79. VoPham T, Hart JE, Laden F, Chiang YY. Emerging trends in geospatial artificial intelligence  
834 (geoAI): potential applications for environmental epidemiology. *Environ Health*. 2018;17.
- 835 80. Mooney SJ, Pejaver V. Big data in public health: terminology, machine learning, and privacy.  
836 In: Fielding JE, Brownson RC, Green LW, editors. *Annual review of public health 2018*. p. 95–112.
- 837 81. Lin Y, Chiang Y-Y, Pan F, Stripelis D, Ambite JL, Eckel SP, et al. Mining public datasets for  
838 modeling intra-city PM2.5 concentrations at a fine spatial resolution. *Proceedings of the 25th ACM*  
839 *SIGSPATIAL international conference on advances in geographic information systems*; Redondo  
840 Beach, CA, USA. 3140013: ACM; 2017. p. 1-10.
- 841 82. Montresor A, Crompton DW, Hall A, Bundy D, Savioli L, Organization WH. *Guidelines for the*  
842 *evaluation of soil-transmitted helminthiasis and schistosomiasis at community level: a guide for*  
843 *managers of control programmes*. Geneva: World Health Organization; 1998.
- 844 83. Burns CJ, Wright M, Pierson JB, Bateson TF, Burstyn I, Goldstein DA, et al. Evaluating  
845 uncertainty to strengthen epidemiologic data for use in human health risk assessments. *Environ*  
846 *Health Perspect*. 2014;122:1160–65.

847

## 848 **Figure Legends**

849

850 **Fig 1.** Study Area: The Mindanao region in The Philippines. Blue dots are the aggregated  
851 survey data at barangay-level.

852 **Fig 2.** Environmental risk factors extraction at pixel-level from cities within barangays

853 **Fig 3.** Posterior estimates and their credible intervals: a) Normalized difference vegetation  
854 index; b) Normalized difference water index; c) Land surface temperature day d) Land surface  
855 temperature night; e) Elevation; f) Nearest distance to water bodies. *Abbreviations:* SSA,  
856 Spatial support of analysis

857 **Fig 4.** Density plots for the risk factors regression coefficients: a) Normalized difference  
858 vegetation index; b) Normalized difference water index; c) Land surface temperature day d)  
859 Land surface temperature night; e) Elevation; f) Nearest distance to water bodies

860 **Fig 5.** Residual plot for the five increasing spatial supports of analysis. The x axis represents  
 861 the fitted prevalence values for the five spatial supports of analysis. The y axis represents the  
 862 residuals calculated by the difference between the observed and predicted prevalence values.

863 **Fig 6.** Proportion of simulated prevalence data that fit the observed maximum prevalence  
 864 value. a) SSA=30m, b) SSA=90m, c) SSA=250m, d) SSA=500m, e) SSA=1 km.  
 865 *Abbreviations:* SSA, Spatial support of analysis.

866 **Fig 7.** Proportion of simulated prevalence data that fit the observed minimum prevalence  
 867 value. a) SSA=30m, b) SSA=90m, c) SSA=250m, d) SSA=500m, e) SSA=1 km.  
 868 *Abbreviations:* SSA, Spatial support of analysis.

869 **Fig 8.** Proportion of simulated prevalence data that fit the observed mean prevalence value. a)  
 870 SSA=30m, b) SSA=90m, c) SSA=250m, d) SSA=500m, e) SSA=1 km. *Abbreviations:* SSA,  
 871 Spatial support of analysis.

## 872 **Tables**

874 **Table 1.** Environmental variables description

Environmental variable	Spatial resolution	Temporal resolution	Data Type	Original coordinate system	Data Source
Elevation	30 m	NA	Raster	EPSG:4326	ASTER GDEM V2 from USGS
NDVI	250 m	2008	Raster	EPSG:4326	MOD13Q1
NDWI	30 m	2008	Raster	EPSG:4326	Landsat 5
	500 m	2008	Raster	EPSG:32651	Landsat 7, 1-year composite
LST	30 m	2008	Raster	EPSG:4326	Landsat 5
	1 km	2008	Raster	EPSG:4326	MOD11A2
NDWB	250 m	2010	Raster	EPSG:32651	Derived from closest facility network using roads, urban areas, river network, and water bodies

875 NDVI: Normalized difference vegetation index; NDWI: Normalized difference water index; LST: Land surface temperature  
 876 day and night; NDWB: Nearest distance to water bodies. USGS: United States Geological Survey.

877

878

879

880

881

882

883

884

885

**Table 2.** Regression coefficient estimates for each risk factor at five descending spatial supports of analysis

Spatial Supports of analysis	Posterior Mean (95% CrI)						
	Intercept	NDVI	NDWI	LSTD	LSTN	E	NDWB
30 m	-6.17 (-6.6, -5.71)	0.21 (-0.05, 0.48)*	-1.06 (-1.3, -0.8)	-0.71 (-1, -0.41)	-0.78 (-1.06, -0.51)	-0.95 (-1.25, -0.65)	-0.28 (-0.43, -0.12)
90 m	-6.16 (-6.61, -5.71)	0.19 (-0.1, 0.48)*	-1.06 (-1.3, -0.8)	-0.71 (-1.01, -0.39)	-0.82 (-1.12, -0.55)	-1.01 (-1.3, -0.71)	-0.29 (-0.45, -0.12)
250 m	-6.08 (-6.48, -5.69)	0.25 (-0.06, 0.54)*	-1.02 (-1.25, -0.76)	-0.71 (-1.02, -0.38)	-0.82 (-1.11, -0.52)	-1.03 (-1.37, -0.71)	-0.31 (-0.5, -0.13)
500 m	-6.02 (-6.43, -5.61)	0.26 (-0.01, 0.53)*	-0.8 (-1.01, -0.57)	-0.65 (-0.97, -0.33)	-0.86 (-1.14, -0.56)	-1.03 (-1.34, -0.71)	-0.38 (-0.57, -0.19)
1 km	-5.49 (-5.76, -5.21)	0.59 (0.3, 0.88)	-0.76 (-1, -0.54)	-0.11 (-0.85, 0.64)*	0.1 (-0.46, 0.68)*	-0.04 (-0.21, 0.14)*	-0.4 (-0.59, -0.2)

Variance of spatial random effect	$\varphi$	DIC
2.5 (1.65, 3.58)	$8.02 \times 10^{-5}$ (-0.004, 0.004)	127.7
2.43 (1.57, 3.5)	$1.2 \times 10^{-4}$ (-0.004, 0.004)	140.5
2.46 (1.6, 3.55)	$9.86 \times 10^{-5}$ (-0.004, 0.004)	86.67
2.35 (1.4, 3.54)	$2.81 \times 10^{-4}$ (-0.004, 0.004)	147.5
2.7 (1.87, 3.7)	$1.65 \times 10^{-5}$ (-0.004, 0.004)	143.7

NDVI: Normalized difference vegetation index; NDWI: Normalized difference water index; LSTD: Land surface temperature day; LSTN: Land surface temperature night; NDWB: Nearest distance to water bodies; CrI: Credible Interval.

\*Non-significant variable in the model

**Table 3.** Credible interval widths (Uncertainty) at five increasing spatial supports of analysis

Spatial Supports of analysis	Credible intervals width (Uncertainty)					
	NDVI	NDWI	LSTD	LSTN	E	NDWB
30 m	0.52	0.50	0.59	0.56	0.59	0.32
90 m	0.57	0.50	0.62	0.56	0.59	0.33
250 m	0.60	0.48	0.64	0.59	0.66	0.36
500 m	0.54	0.44	0.64	0.58	0.64	0.38
1 km	0.58	0.46	<b>1.49</b>	<b>1.14</b>	0.34	0.39

NDVI: Normalized difference vegetation index; NDWI: Normalized difference water index; LSTD: Land surface temperature day; LSTN: Land surface temperature night; NDWB: Nearest distance to water bodies.

High uncertainty values are present in bold

**Table 4.** Correlation and determination coefficients between the spatial supports of analysis (SSAs) and environmental covariates estimates and uncertainties.

	Estimates			Uncertainties		
	Correlation coefficient	Determination coefficient (R <sup>2</sup> )	p-value	Correlation coefficient	Determination coefficient (R <sup>2</sup> )	p-value
<b>NDVI</b>	0.94	0.85	0.02	0.32	-0.2	0.6
<b>NDWI</b>	0.93	0.81	0.02	-0.3	-0.2	0.6
<b>LSTD</b>	0.92	0.8	0.03	0.91	0.76	0.03
<b>LSTN</b>	0.86	0.6	0.06	0.9	0.75	0.04
<b>E</b>	0.86	0.65	0.06	-0.78	0.48	0.12
<b>NDWB</b>	-0.94	0.86	0.02	0.91	0.76	0.03
<b>Variance</b>	0.64	0.21	0.25	-0.54	0.06	0.35

**Table 5.** Resulting ppp-values for the test statistics: maximum (8.5%), minimum (0.33%) and mean (0.5%) prevalence values at five increasing SSA.

Spatial Supports of analysis	ppp-value (maximum)	ppp-value (minimum)	ppp-value (mean)
30 m	0.66	0.87	0.66
90 m	0.67	0.86	0.66
250 m	0.66	0.88	0.63
500 m	0.66	0.87	0.67
1 km	0.64	0.93	0.59

ppp-value: posterior predictive *p*-value