



University of Pennsylvania
ScholarlyCommons

Statistics Papers

Wharton Faculty Research

9-11-2014

Mapping the Spatial Distribution of a Disease-Transmitting Insect in the Presence of Surveillance Error and Missing Data

Andrew E. Hong

Corentin M. Barbu

Dylan S. Small
University of Pennsylvania

Michael Z. Levy
University of Pennsylvania

Follow this and additional works at: https://repository.upenn.edu/statistics_papers

 Part of the [Business Analytics Commons](#), [Management Sciences and Quantitative Methods Commons](#), and the [Statistics and Probability Commons](#)

Recommended Citation

Hong, A. E., Barbu, C. M., Small, D. S., & Levy, M. Z. (2014). Mapping the Spatial Distribution of a Disease-Transmitting Insect in the Presence of Surveillance Error and Missing Data. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 178 (3), 641-658. <http://dx.doi.org/10.1111/rssa.12077>

This paper is posted at ScholarlyCommons. https://repository.upenn.edu/statistics_papers/636
For more information, please contact repository@pobox.upenn.edu.

Mapping the Spatial Distribution of a Disease-Transmitting Insect in the Presence of Surveillance Error and Missing Data

Abstract

Maps of the distribution of epidemiological data often ignore surveillance error or possible correlations between missing information and outcomes. We analyse presence–absence data at the household level (12050 points) of a disease-carrying insect in Mariano Melgar, Peru, collected as part of the Arequipan Ministry of Health's efforts to control Chagas disease. We construct a Bayesian hierarchical model to locate regions that are vulnerable to under-reporting due to surveillance error, accounting for variability in participation due to infestation status. The spatial correlation in the data allows us to identify relative inspector sensitivity and to elucidate the relationship between participation and infestation. We show that naive estimates of prevalence would be biased by surveillance error and missingness at random assumptions. We validate our results through simulations and observe how randomized inspector assignments may improve prevalence estimates. Our results suggests that bias due to imperfect observations and missingness at random can be assessed and corrected in prevalence estimates of spatially auto-correlated binary variables.

Keywords

Bayesian hierarchical modelling, spatial analysis, statistical epidemiology, surveillance error

Disciplines

Business | Business Analytics | Management Sciences and Quantitative Methods | Statistics and Probability

Mapping the Spatial Distribution of a Disease Transmitting Insect in the Presence of Surveillance Error and Missing Data

Andrew E. Hong¹, Corentin M. Barbu², Dylan S. Small¹, Michael Z.
Levy², and The Chagas Disease Working Group in Arequipa

¹Department of Statistics, The Wharton School, University of
Pennsylvania, 400 Jon M. Huntsman Hall, 3730 Walnut Street,
Philadelphia, PA 19104, USA

²Department of Biostatistics and Epidemiology, Center for Clinical
Epidemiology and Biostatistics, University of Pennsylvania, 800
Blockley Hall, Philadelphia, PA 19104, USA

March 26, 2014

Abstract

Maps of the distribution of epidemiological data often ignore surveillance error or possible correlations between missing information and outcomes. We analyse presence-absence data at the household level (12,050 points) of a disease carrying insect in Mariano Melgar, Peru, collected as part of the Arequipan Ministry of Health's efforts to control Chagas disease. We construct a Bayesian hierarchical model to locate regions vulnerable to under-reporting due to surveillance error, accounting for variability in participation due to infestation status. The spatial correlation in the data allows us to identify relative inspector sensitivity and elucidate the relation between participation and infestation. We show that naïve estimates of prevalence would be biased by surveillance error and missing at random (MAR) assumptions. We validate our results through simulations and observe how randomized inspector assignments may improve prevalence estimates. Our results suggests that bias due to imperfect observations and MAR can be assessed and corrected in prevalence estimates of spatially autocorrelated binary variables.

1 Introduction

The increasing risk of vector-borne disease epidemics has accompanied the rise of urban environments in the developing part of the world. The use of spatial analysis in public health campaigns for disease control is documented in Dengue ([1, 2, 3]), Malaria ([4, 5, 6]), and Chagas disease ([7]). Chagas disease is a tropical parasitic disease, affecting millions in Central and South America. The disease agent is *Trypanosoma cruzi*, a parasite that is transmitted by the *Triatoma infestans* insect vector. Policy for Chagas disease control has focused on the elimination of this vector [8]. While initiatives to control *T. Infestans* have been active for decades [9], the insect is a continually re-emergent threat in Peru [10]. Because of the strain on public resources created by these recurring epidemics and the risk of emergence of insecticide resistance due to repeated treatment [11, 12], there is an interest in applying statistical methods to guide the application of insecticide to urban areas [13, 14].

This study was done in coordination with the efforts of the Peruvian Ministry of Health to control an epidemic of *T. cruzi* infections in the city of Arequipa, Peru, [10]. Insecticide treatment was preceded by a household level survey, which identified household infestations of *T. infestans*. The results of the survey, conducted in the district of Mariano Melgar, are shown in Figure 1. At the time of the survey, policy was to prioritize the treatment of households in district localities where the rate of household infestations exceeds ten percent. Ecological surveys for observing presence-absence are conducted by human inspectors and are subject to under-reporting of presence. Inspectors are heterogeneous, differing in their ability to identify infestations. While previous work on triatomine infestation detection has used spatial techniques, acknowledging the imperfect inspection process, this work has not accounted for heterogeneity in the inspectors' skills [14]. Another concern of policymakers in this survey is the large proportion of missing information (34% of the records). Correlation between missingness and the studied outcome have been shown to lead to serious bias in clinical trials, raising serious concerns of the validity of research findings [15]. In the case of spatially autocorrelated outcomes, Bayesian hierarchical models have been used extensively to obtain point estimates at missing locations under the assumption that missingness is uncorrelated with the outcome [16, 17, 18].

Here, we propose to adjust the risk mapping of triatomine infestations for surveillance error, caused by the lack of sensitivity on the part of inspectors and missing not at random, caused by different inclination to par-

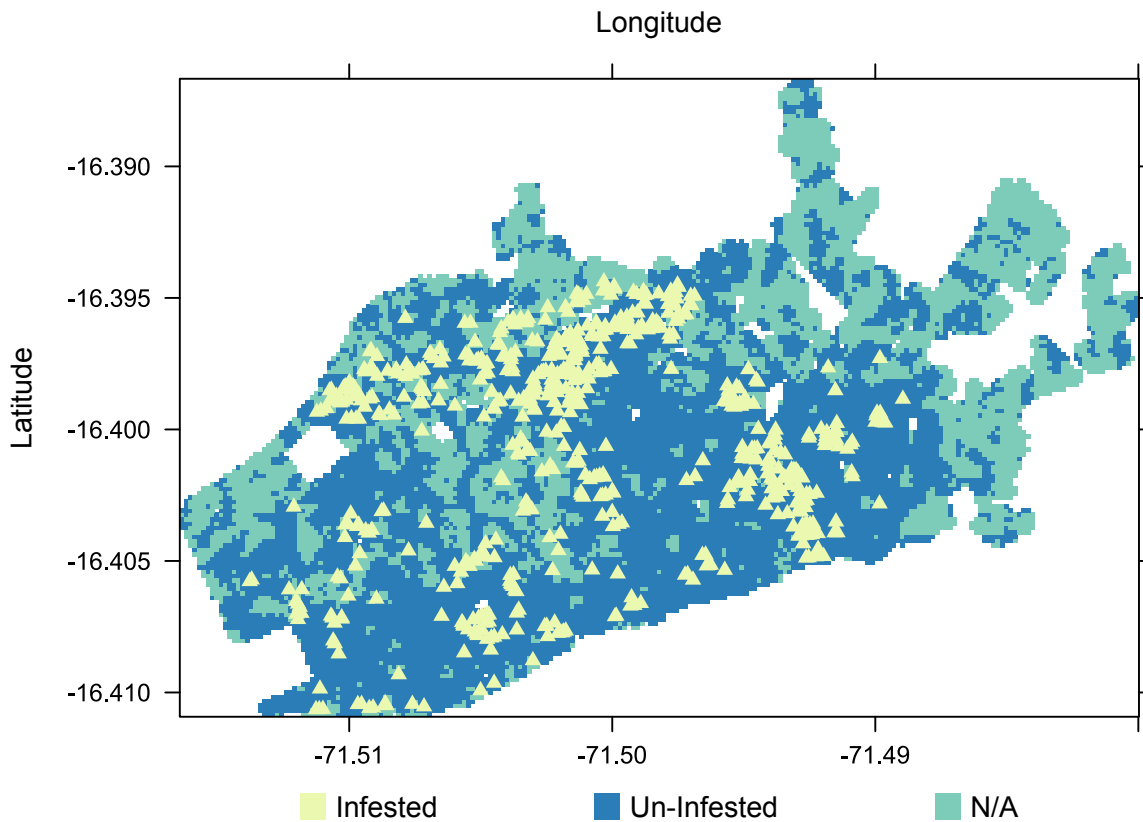


Figure 1: Transect data from the 2011 *T. infestans* survey in Mariano Melgar, Arequipa, Peru that was collected by the Ministry of Health. This area is 3838 meters by 2664 meters and contains 12,050 total households. This survey identified 608 positive households and contains 4,098 non-participating households.

ticipate in the surveys depending on the infestation status. Explicitly, we construct a Bayesian hierarchical model that jointly assess 1) the probability that households participate depending on their infestation status 2) the individual sensitivities of the inspectors and 3) the prevalence of infestation accounting for the former. We then discuss the findings of our model in Mariano Melgar that guided the Ministry of Health’s 2011 campaign in Arequipa.

2 Inspection & Data Collection

The *T. infestans* survey was conducted by the Ministry of Health in Mariano Melgar, a district of 12,050 households. We mapped the locations of these households by determining their relative position to city blocks and comparing field maps of these blocks to satellite images from Google EarthTM, [19]. Inspectors requested the participation of residents before searching households for *T. infestans*. The outcome of each inspection was either the successful collection of insect samples, the failure to locate samples, or non-participation. Each entry of the data consists of: a pair of coordinates denoting the location of the household, a presence-absence status, and the identifier or labelling of the inspector. 4,098 households (around 34 %) opted not to participate in the survey. In this study, missingness appears to aggregate spatially in regions with lower rates of infestation and is therefore missing not at random (MNAR). We model the relationship between the true infestation and the point pattern of missingness to analyze this claim in Section 3.2.

Separate data identifying the sensitivity of the 40 inspectors involved in this study was unavailable. However, validated data from previous treatment campaigns suggested the general sensitivity of human inspectors, which informed our prior specification. We rely on knowledge of the coordinate locations of the households to infer the distribution of the *T. infestans* infestation and inspectors assignments. The household assignments of four inspectors is shown in Figure 2. During the survey, inspectors were assigned to households, based on staffing constraints, which resulted in subgroups of inspectors inspecting an entire region. Because of these aggregated assignments, the surveillance error in this study is spatially correlated. This confounding between infestation distribution and inspector sensitivity through geographic location potentially biases the estimation of insect presence. For a geographic region, it becomes difficult to disassociate the severity of the infestation apart from the sensitivity of the inspectors. For this reason, we

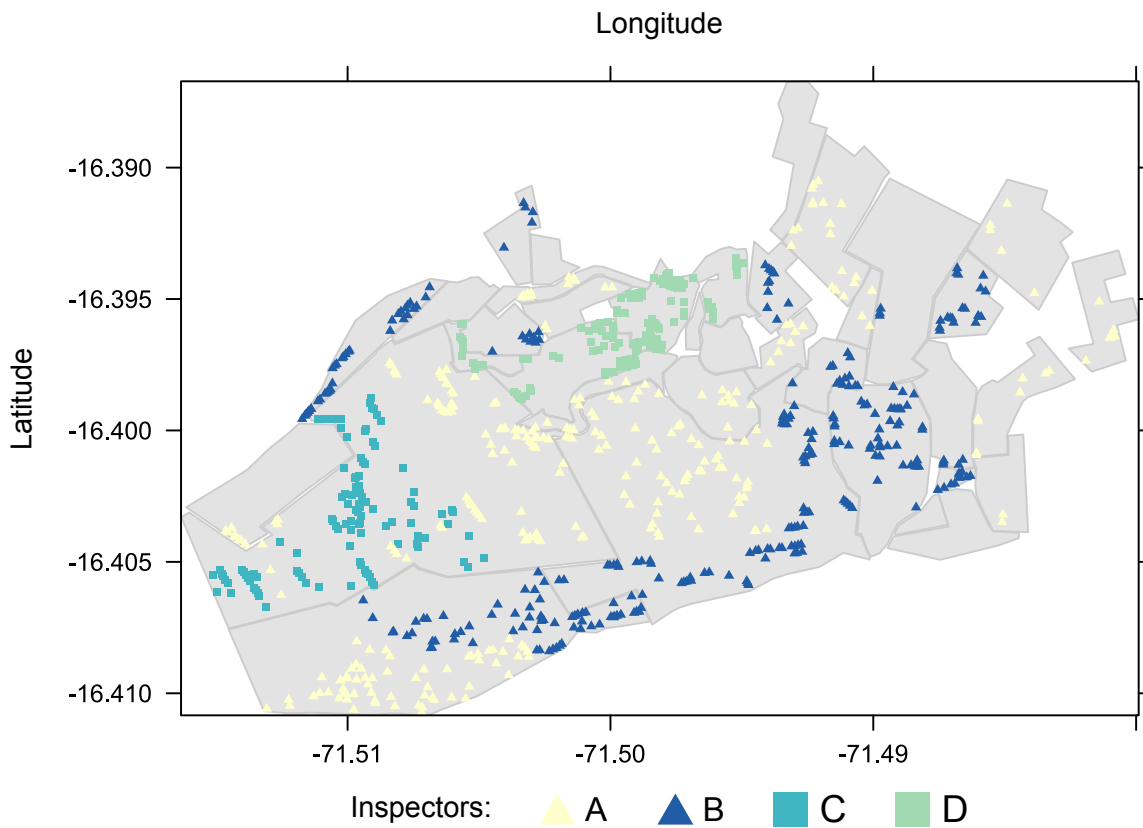


Figure 2: Examples of the household assignments of four inspectors across the 35 localities of Mariano Melgar. The number and distribution of assigned households inspected varies by individual. The assignments of inspectors A (301 total households) and B (291) span multiple localities across the district, whereas the assignments of inspectors C (107) and D (137) are highly localized.

study in Section 5 how prevalence estimates may be subject to confounding by inspectors assignment.

3 Model Specification

The primary interest in this study is to infer the true household infestation status, which we modelled as a binary outcome: infested or un-infested. In the context of our application, this quantity is not directly observable through the observation of inspectors, who are subject to surveillance error. We model the true binary infestation status using a generalized linear model with the probit link function. A complete diagram of the hierarchical components is shown in Figure 3. Each of the 12,050 households is included in the model and indexed by i . The i th household is located at a pair of points denoting easting and northing in the Universal Transverse Mercator projection coordinate system. For household i , the variable y_i is the true infestation status, where $\{y_i = 1\}$ indicates that household i is infested. We model the probability that the i -th household is infested by the following,

$$\mathbb{P}(y_i = 1|u_i, t) = \Phi(u_i + t) \tag{1}$$

where Φ is the cumulative distribution function of the standard Gaussian. u_i is a continuous, household level effect, capturing how at risk the i -th household is for being infested, and t is an intercept term for the entire district. The covariance structure for $u = [u_i]_{i=1}^n$ among the households is defined by their geographic locations and will ensure that the infestation statuses are spatially correlated. The usual convention is to place a diffuse $N(0, \sigma^2)$ prior on t . Similar approaches for spatial modelling of categorical data may be found in [20] and [21] for other public health settings.

3.1 Spatial Effect

We used a conditionally auto-regressive Gaussian model for u to capture the spatial similarity of the infestation. The model for u , popularized by [22], is a centered Gaussian with a precision matrix Λ . The entries of the model precision matrix are based on the pairwise euclidean distance between households $d_{i,j}$, a scaling parameter k_u , and a threshold T .

$$\Lambda_{i,j} = \begin{cases} k_u \sum_{\{k:d_{i,k}<T\}} d_{i,k}^{-1}, & \text{if } i = j \\ -k_u d_{i,j}^{-1} \mathbb{1}_{d_{i,j} \leq T}, & \text{if } i \neq j \end{cases} \tag{2}$$

The effect of the scaling and threshold parameters is most evident on the marginal distributions of u_i conditional on the rest of the households u_{-i} , [23]:

$$u_i|u_{-i} \sim N\left(\frac{\sum_{\{j:d_{i,j}<T\}} d_{i,j}^{-1} u_j}{\sum_{\{j:d_{i,j}<T\}} d_{i,j}^{-1}}, \frac{1}{k_u \sum_{\{j:d_{i,j}<T\}} d_{i,j}^{-1}}\right) \quad (3)$$

For the conditional marginal, $u_i|u_{-i}$ is centered at a weighted sum of neighbouring values within the threshold radius. Households, whose distance to i exceeds T , have no effect on the conditional distribution of u_i . Within the threshold radius T , households closer to i are given more weight proportional to their inverse distance. The scaling parameter k_u determines the variation of u_i around this center.

Previously, we had analysed a neighbouring district in Arequipa and determined that the correlation in infestation statuses is negligible for households separated by 50 meters or more, [24]. We fix our threshold T at 50 meters. For the prior on k_u we follow the usual practice of placing a conjugate $\exp(\lambda)$ prior on k_u [25].

3.2 Missing Data

We treat the missingness point pattern, which we denote as $\mathbb{1}_{NA}$, as a series of Bernoulli outcomes, depending on the true infestation status of the household. The probability of household participation is modelled separately for infested $\pi_{NA}^{(1)}$ and un-infested households $\pi_{NA}^{(0)}$ allowing for differential participation according to the infestation. Further, due to the marked differences of socio-economic status between localities, we allow this relationship to vary across localities:

$$\pi_{NAj}^{(1)} = \mathbb{P}(\mathbb{1}_{NAi} = 0 | y_i = 1) \quad (4)$$

$$\pi_{NAj}^{(0)} = \mathbb{P}(\mathbb{1}_{NAi} = 0 | y_i = 0) \quad (5)$$

where the household indexed by i is located in the locality j . We used identical beta distribution, $B(p, q)$, priors for both parameters in each of the localities.

3.3 Surveillance Process

To account for the human error in surveillance, we model the sensitivity of each inspector as the probability, $\beta \in [0, 1]$, that an inspector locates

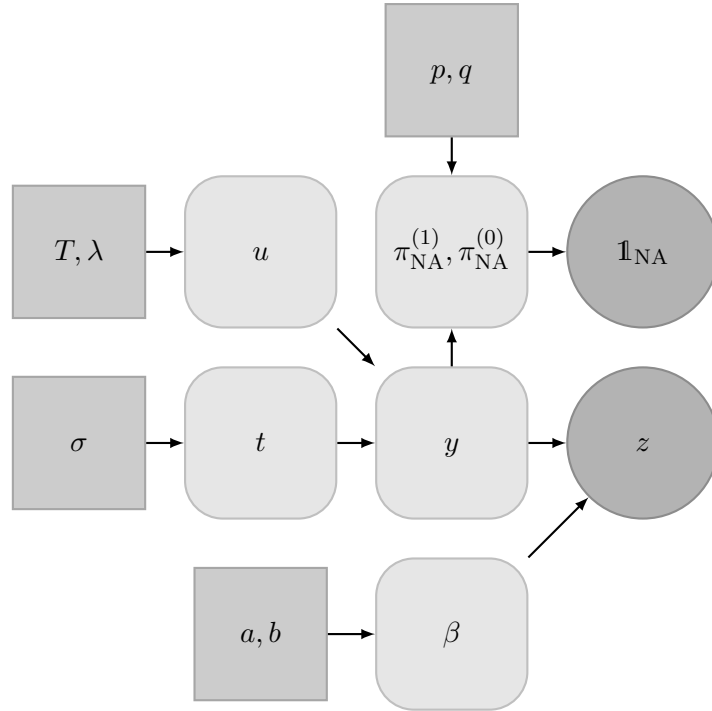


Figure 3: Model diagram. The variables z , the reported infestation status, and $\mathbb{1}_{NA}$, the non-response indicator, are the observed data. z is generated by y , the true infestation status, and β , the inspector sensitivity. $\mathbb{1}_{NA}$ is generated by y , the true infestation status, and $(\pi_{NA}^{(0)}, \pi_{NA}^{(1)})$ the probability of participation depending on the true infestation status. The main parameter of interest is the infestation status of the households, y , which is spatially correlated through the random field u .

T. infestans in the household, when the insect is present. Then, the reported outcome: infested or un-infested follows a Bernoulli distribution in observed households. If $\beta(i)$ is the sensitivity of the inspector that inspected household i , the distribution of the reported outcome is as follows,

$$\mathbb{P}(z_i = 1 | y_i, \beta(i), \mathbf{1}_{\text{NA}i}) = \begin{cases} \beta(i)y_i, & \text{if } \mathbf{1}_{\text{NA}i} = 0 \end{cases} \quad (6)$$

The sensitivity, $\beta(i)$, is therefore only relevant if the true infestation status of the household, y_i , is positive. This surveillance model is the individual inspector sensitivity model, where each inspector in the study has his or her own sensitivity parameter, β_j . In contrast to this individual inspector model, a simpler model is the group inspector model, where the sensitivity is identically β for all of the households. We use a beta distribution, $B(a, b)$, as the prior for the sensitivity parameters either for each inspector separately or for the common sensitivity of all the inspectors in the simpler model.

4 Results for the 2011 Mariano Melgar Survey

Infestation by *T. infestans* in Mariano Melgar, Arequipa was strongly clustered in space. Based on the raw surveys (proportion of infested among the surveyed households), only four localities fit the 10% prevalence criterion for inclusion in blanket insecticide treatment, which is the uniform application of insecticide to all households in the locality. Our role in this study was to apply the model proposed in Section 3 to adjust these estimates and possibly identify additional localities at risk for major *T. infestans* infestations.

4.1 Full model results

To perform the analysis on the Mariano Melgar survey, we placed priors on the following parameters: the intercept of the linear model, t ; the precision parameter of the Gaussian spatial effect, k_u ; and the inspector sensitivities, $\{\beta_j\}_j$. For the first two parameters, we used a diffuse Gaussian, $N(0, \sigma^{-2} = 1e-8)$, prior and a diffuse exponential, $\exp(\lambda = 1e-4)$, prior. Our estimates of the presence absence values are dependent on the specification of the inspector sensitivity priors. The sensitivity of our analysis to this prior is displayed in Table 4. Cross-sectional pre- and post-treatment data from previous spraying campaigns showed that human inspectors are accurate roughly 76 percent of the time. After consulting local experts, we agreed upon the use of a beta, $B(6.5, 2)$, prior for inspectors' sensitivities. We used a relatively weak beta, $B(7, 3)$, prior for the missingness probabilities

to reflect the overall rate of missing households in the data (around 0.34). We found that the choice of prior made little change to the Monte Carlo estimates for the Mariano Melgar data.

We implemented the model using a Gibbs sampler, which is outlined in Appendix A, and estimated the posterior probability of infestation, $\hat{\mathbb{P}}(y_i = 1|z)$, for each of the households. Averaging these household level estimates by locality, we produced the locality-wide infestation estimates shown in Table 1. With our informative prior for inspector sensitivity, we found that the locality estimates for two additional localities, 11 and 37, exceeded the 10 percent mark. We mapped the infestation estimates at the block level in Figure 4 in order to provide guidance on insecticide application in areas where blanket treatment (insecticide application to all participating households) is not warranted.

By introducing dependence between the missingness point pattern and the infestation, data is not missing at random. For the Mariano Melgar survey, we allowed this dependence to vary across localities. This model inferred that survey participation rates were higher for infested households compared to un-infested households, consistently across localities. Across almost all localities, we found that the estimated participation rates, displayed in Table 3, were higher for infested households (with the exception of locality number 9). Because the model linked under-participation to lower infestation rates, it is likely that estimates of prevalence made under the MAR assumption would overestimate of the infestation.

Figure 5 displays the posterior distributions of the least sensitive, the 10th, 20th, 30th, and most sensitive inspectors, ranked by posterior mean. While these posterior distributions vary from the prior, the group average of all the inspectors' posterior means was 0.75601, which was close to the prior mean of 0.7647. Similarly, when using $B(1/2, 1/2)$ and $B(1, 1)$ priors of mean 0.5 we found that the group averages were 0.5663 and 0.5507 respectively. While the overall levels of estimated infestation were sensitive to the prior, we found that the rankings of inspectors remained consistent across prior specifications, see Table 4. This consistency suggests that there is information present in the data to identify the relative sensitivity of inspectors.

4.2 Model Comparison

We compared the Bernoulli inspection model in Section 3.3 to a more traditional regression model. We model the reported infestation outcome (binary), y , as,

$$y_{ij} = x_i + \beta_j \mathbb{1}_{NA_i} + t \tag{7}$$

where i denotes the location and j denotes the inspector. Because the outcome is usually taken as observed unambiguously, we make the comparison between the two approaches based on their infestation estimate for the missing households in the study. We treat the non-participating households, identically to the inspected households in the study, except for the fact that these households have a fixed inspector effect equal to zero (the effect of an average inspector with an effect equal to the prior mean). We used the same models described in 3 for the spatial effect, x , and the intercept, t . In contrast to Section 3.3, the inspector effect, β , is a continuous, unbounded variable. In practice, we placed diffuse, centered Gaussian priors on these parameters, unlike the beta priors used for the Bernoulli model.

For the Mariano Melgar study, the estimates for the non-participating households using this more traditional regression model were similar to the estimates produced by our model. The estimates are particularly similar between a more traditional approach and our model when households are assumed to participate at random, demonstrated in Table 1. This similarity is expected as non-participating households are treated as locations to be interpolated and not used to inform the fit of the model. However, the traditional regression approach does not estimate the negative households probability to be false negative depending on their inspector.

Table 1: Estimated Prevalence of *T. Infestans* in localities of the district of Mariano Melgar

	# Units	# Infest.	# NA	Prop. of Infected Participants	Infest. Estimate (MAR)	Infest. Estimate (MNAR)
1	294	0	173	0	0.0041	0.0025
2	2605	142	998	0.0884	0.1092	0.0821
3	271	5	125	0.0342	0.0755	0.0559
4	170	0	60	0	0.0056	0.0051
5	82	0	49	0	0.0133	0.0032
6	73	0	45	0	0.0101	0.0043
7	82	0	51	0	0.0994	0.0402
8	37	0	12	0	0.0103	0.0076
9	108	16	1	0.1495	0.2016	0.1916
10	147	17	77	0.2429	0.3077	0.2163
11	132	8	38	0.0851	0.1200	0.1104
12	113	18	56	0.3158	0.3784	0.2880
13	604	145	50	0.2617	0.3548	0.3378
14	147	0	60	0	0.0147	0.0154
15	273	16	78	0.0821	0.1027	0.0885
16	134	3	57	0.0390	0.0339	0.0310
17	113	0	53	0	0.0260	0.0134
18	180	0	101	0	0.0023	0.0009
19	169	0	100	0	0.0005	0.0005
21	374	0	203	0	0.0003	0.0003
22	108	0	64	0	0.0014	0.0019
23	176	1	59	0.0085	0.0198	0.0217
24	225	1	101	0.0081	0.0125	0.0089
25	166	0	83	0	0.0038	0.0008
26	226	0	147	0	0.0055	0.0012
28	143	0	82	0	0.0072	0.0000
30	213	0	122	0	0.0018	0.0010
31	106	0	62	0	0.0040	0.0022
32	33	0	12	0	0.0031	0.0016
33	82	0	57	0	0.0074	0.0010
34	50	0	22	0	0.0000	0.0000
35	160	0	87	0	0.0087	0.0026
36	534	20	96	0.0457	0.0572	0.0569
37	2022	137	373	0.0831	0.1086	0.1000
38	1698	79	344	0.0583	0.0770	0.0706

Displayed above are the total number of households, number of positively identified households for infestations, number of non-participating households, and average probability of infestation by locality. We used these estimates to identify two additional at-risk localities: 11 and 37 that were later treated with insecticide in the Spring of 2011.

Table 2: Comparison Infestation Estimates for Missing Households between the Bernoulli Inspection Error Model and Traditional Regression

	Estimate (NA) Traditional Regression	Estimate (NA) MAR	Estimate (NA) MNAR
1	0.0023	0.0055	0.0018
2	0.1020	0.0995	0.0223
3	0.0986	0.0800	0.0358
4	0.0171	0.0109	0.0101
5	0.0219	0.0164	0.0019
6	0.0205	0.0125	0.0043
7	0.1228	0.1428	0.0542
8	0.0217	0.0125	0.0117
9	0.0500	0.0300	0.1200
10	0.2510	0.2810	0.0930
11	0.0945	0.1019	0.0553
12	0.3388	0.3476	0.1638
13	0.2224	0.2880	0.0962
14	0.0334	0.0267	0.0226
15	0.0829	0.0816	0.0294
16	0.0055	0.0132	0.0055
17	0.0461	0.0412	0.0174
18	0.0003	0.0035	0.0008
19	0.0021	0.0008	0.0005
21	0.0024	0.0005	0.0004
22	0.0066	0.0021	0.0018
23	0.0345	0.0297	0.0319
24	0.0141	0.0139	0.0061
25	0.0032	0.0058	0.0008
26	0.0026	0.0076	0.0011
28	0.0025	0.0107	0.0000
30	0.0013	0.0026	0.0010
31	0.0076	0.0062	0.0028
32	0.0025	0.0025	0.0000
33	0.0041	0.0086	0.0004
34	0.0000	0.0000	0.0000
35	0.0137	0.0148	0.0037
36	0.0540	0.0503	0.0422
37	0.0869	0.0869	0.0294
38	0.0806	0.0687	0.0325

This table compares the results produced by our model to a more traditional regression model, where inspectors are treated as fixed regression effects. We average the household estimates for the non-participating households only by locality. The traditional regression estimates were quite similar to the estimates produced by the Bernoulli model, when data assumed to be missing at random.

Table 3: Estimated participation rates in entomological surveys among residents of infested and un-infested households by locality in the district of Mariano Melgar

	π_{NA} for Infested Units	π_{NA} for Un-Infested Units
1	0.7903	0.4267
2	0.8978	0.5934
3	0.7681	0.5364
4	0.7814	0.6510
5	0.7945	0.4509
6	0.7827	0.4252
7	0.6753	0.4404
8	0.7917	0.7075
9	0.9255	0.9716
10	0.7953	0.4339
11	0.8355	0.6959
12	0.7522	0.4578
13	0.9692	0.8832
14	0.7250	0.6029
15	0.8675	0.7000
16	0.8288	0.5841
17	0.7588	0.5552
18	0.8047	0.4538
19	0.7962	0.4269
21	0.8060	0.4692
22	0.7856	0.4389
23	0.7271	0.6790
24	0.7889	0.5621
25	0.7864	0.5188
26	0.7953	0.3686
28	0.7952	0.4562
30	0.8085	0.4446
31	0.7875	0.4427
32	0.7994	0.6623
33	0.8025	0.3591
34	0.7955	0.6110
35	0.7840	0.4783
36	0.8546	0.8170
37	0.9409	0.8015
38	0.9001	0.7872

By introducing dependence between the missingness point pattern and the infestation, we move to a more realistic setting, where data is no longer missing at random. For the Mariano Melgar survey, we allowed this dependence to vary across localities. This model inferred that the participation rates for the survey were higher for infested households compared to un-infested households, consistently across localities.

Table 4: Pairwise Correlations between Estimated Inspector Rankings in the Mariano Melgar Survey against Various Prior Specifications for Surveillance Error

Correlations between Estimated Rankings				
Prior	B(6.5, 2)	B(5, 5)	B(1, 1)	B(1/2, 1/2)
B(6.5, 2)		0.9600	0.9392	0.7356
B(5, 5)	0.9600		0.9765	0.7921
B(1, 1)	0.9392	0.9765		0.7471
B(1/2, 1/2)	0.7356	0.7921	0.7471	

Different prior specifications for the sensitivity parameters produce different estimates of each inspector’s sensitivity when analyzing the Mariano Melgar survey. However, when ranking inspectors based on their estimated sensitivity using these priors, we found that the rankings were fairly consistent.

5 Simulation Study: Household Assignments & Confounding

Confounding, in the context of this study, is the systemic correlation across households between the risk of infestation and the surveillance error of the surveying inspector. Such confounding may potentially bias estimates of inspector sensitivity and infestation in our model. Confounding may be avoided if inspectors were randomly assigned to locations. However, randomized assignment is difficult to implement in practice as it is costly for each inspector to inspect geographically dispersed locations. Inspectors in Mariano Melgar were not randomly assigned to locations; instead, the locations inspected by an inspector tended to be geographically aggregated, see Figure 2. Barring careful assignment of inspectors, it is difficult to ascertain if or to what degree a particular study is affected by confounding. We conduct a simulation study to understand how the inspector assignment in Mariano Melgar affects the estimation of the infestation and inspector sensitivities by comparison to estimation when inspectors are randomly assigned.

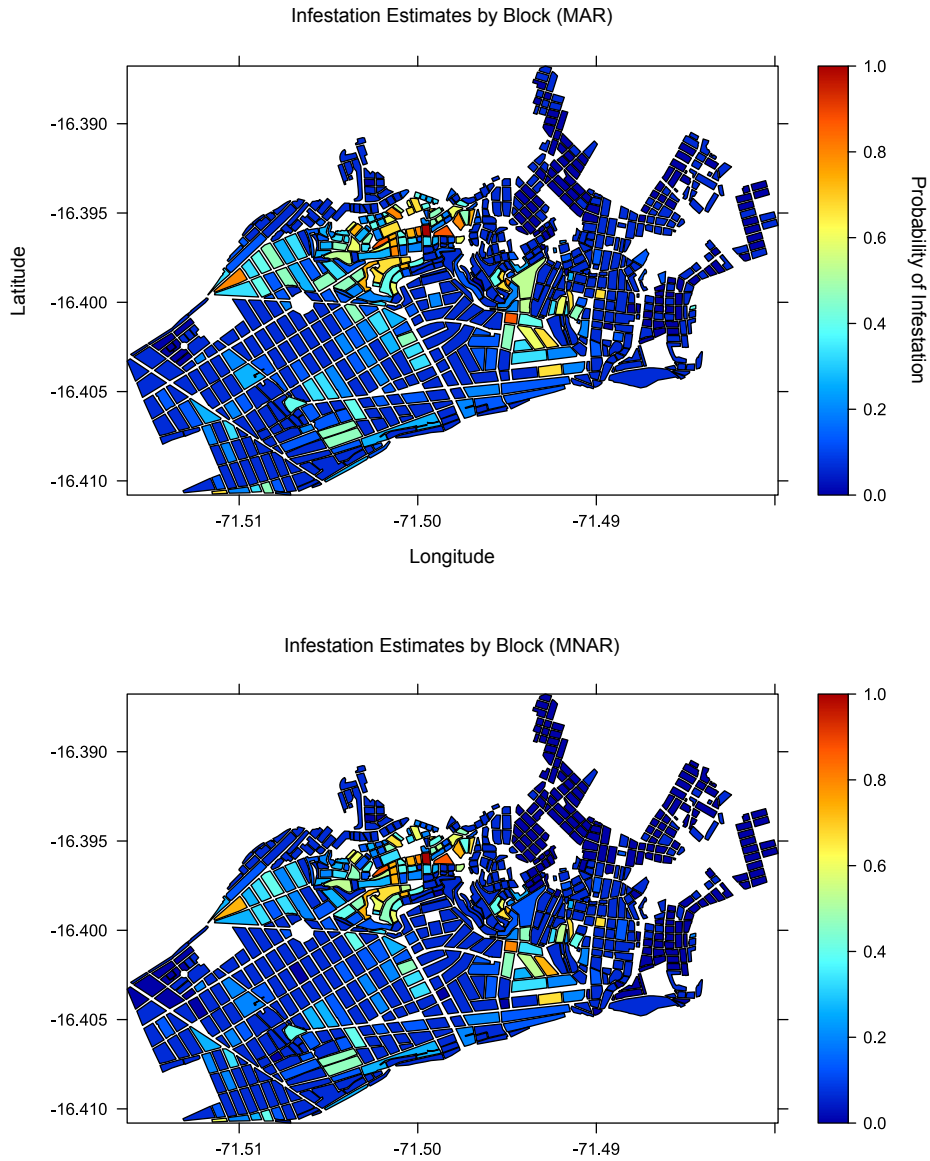


Figure 4: Map of the estimated prevalence of infestation across the city blocks of Mariano Melgar, Arequipa, Peru, prior to insecticide treatment in 2011. The top map displays the estimated infestation prevalence of each household averaged by city block, when survey participation was assumed to be missing at random (MAR). The bottom map displays the estimates produced by our model under the assumption that data was missing not at random (MNAR).

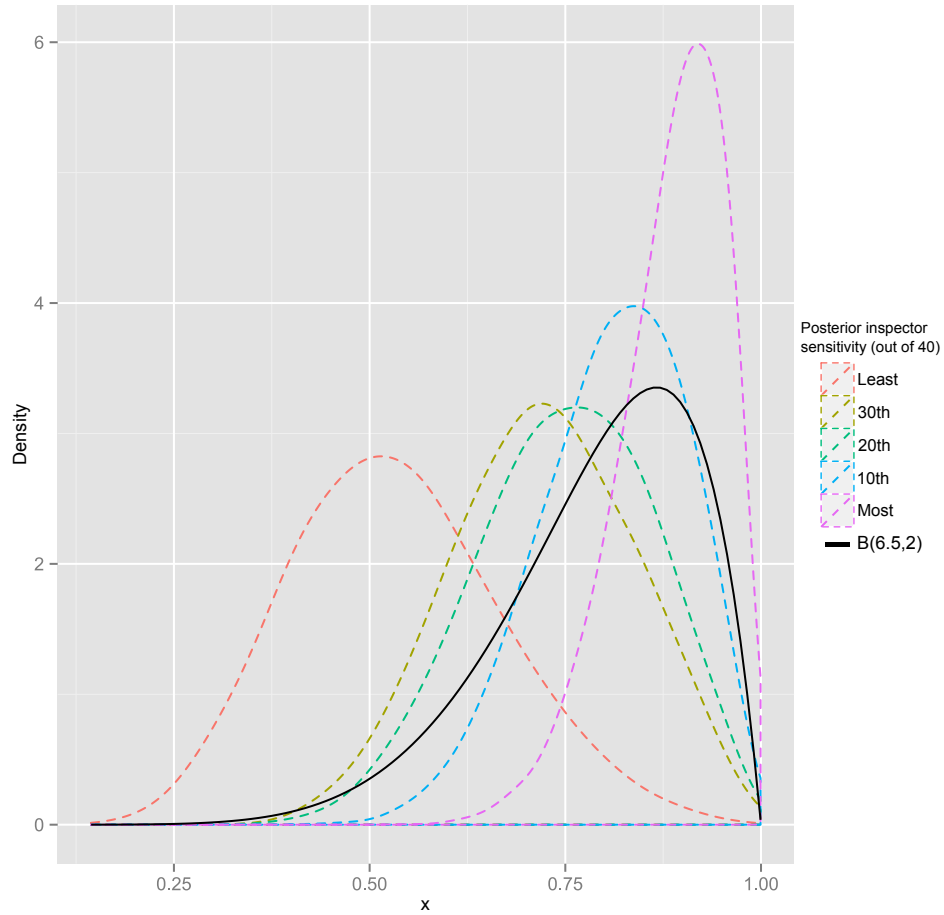


Figure 5: Posterior distributions (dotted lines) of the inspector sensitivity parameters, β . Displayed here are the distributions of the least sensitive, the 10th, 20th, 30th, and most sensitive inspectors from the Mariano Melgar survey, ranked according to their mean sensitivity. These posteriors deviate from the $B(6.5, 2)$ prior for inspector sensitivities (solid line) that was used for the analysis.

5.1 Methodology

Forty inspectors collected the data in the Mariano Melgar survey, each assigned to a set of households in the survey. To create a randomized assignment, we simulate data using the same number of inspectors and totals of households inspected by each inspector, but select uniformly at random the locations of each inspector’s assigned households. The interest of these simulations is the relationship between the spatial distribution of inspectors and our ability to accurately infer the infestation. For these simulations, we do not simulate the missing data point pattern using the model described in Section 3.2 for simplicity.

We simulate a mock infestation using the posterior for true presence, y , from the survey data. To simplify the effect of prior specification, we simulate all inspector sensitivities from a common $B(6.5, 2)$ prior. We then generate two distinct data sets for the observed infestation using the Mariano Melgar survey assignment and the randomized assignment. Finally, we estimate for each data set the household level infestation probability and the inspector sensitivity using the Gibbs sampler (see Appendix A).

5.2 Factors Influencing Estimation Performance

In addition to the effect of the inspector assignments, we are also interested in the effect of inspector sensitivity priors and inspector sensitivity models on the performance of our estimation. As a gold standard, usage of the $B(6.5, 2)$ generating prior should result in the lowest estimation error. We contrast the usage of this prior with the usage of the uniform $B(1, 1)$ prior and the centered $B(5, 5)$ prior. Results using the $B(1, 1)$ demonstrate how effectively the inspector sensitivities may be learned from the data. Usage of the $B(5, 5)$ prior provides insight into the sensitivity of our model to a strongly misspecified prior.

5.3 Estimators & Measures

Each round of the simulation produces a simulated observed presence-absence, z_{Sim} , given the inspector assignment. For each simulated data set, we then estimate the infestation and inspector accuracies for every combination of prior and model factors previously detailed.

The estimate for the infestation is the posterior probability of infestation for each household. The measure of accuracy for the infestation estimates is the squared distance between the simulated infestation data, y_{Sim} and its estimate, \hat{y} . As the former is a binary outcome and the latter is a probability

forecast, the root-squared distance, $\|y_{\text{Sim}} - \hat{y}\|_2$, is the Brier score, [26]. For the inspector sensitivities, the estimate is the posterior expectation $\hat{\beta}$. We again use the squared distance between this estimate and the generating parameter, $\|\hat{\beta} - \beta_{\text{Sim}}\|_2$, to measure the accuracy of our inspector sensitivity estimates.

5.4 Results of the simulation studies

The results of the fifty-run simulation study are summarized for infestation estimation in Tables 5 and inspector sensitivity estimation in Table 6. For every choice of prior, we attained better estimation of the infestation when inspectors were randomly assigned to households compared to when inspectors followed the Mariano Melgar survey assignment and missingness was ignored. However, the differences in Brier score between assignment type were significant but not extreme. In the worst case, when a $B(5, 5)$ prior is placed on the inspector sensitivities, we found that the mean Brier score was only 1.3% larger for the Mariano Melgar assignment compared to randomized assignments. The randomized assignment of inspectors is well-guarded against confounding, where the assignment of inspectors to households is highly spatially correlated, but the performance of our model under the survey assignments was not significantly worse. These results demonstrate that the inspector assignment used to conduct the Mariano Melgar survey is not prone to excessive confounding error and lend credence to our findings in Section 4.

These simulations also affirm the sensitivity of the surveillance error estimation to prior specification. For the actual survey assignment of inspectors, the mean estimation error for the inspector sensitivities using the misspecified $B(5, 5)$ prior was 79% larger than the mean using the correct $B(6.5, 2)$ prior. Similarly, the mean estimation error using the weak $B(1, 1)$ prior was 14% larger than the mean error using the $B(6.5, 2)$ prior.

We conclude based on the significance of these results that there is insufficient information in the data to infer the absolute sensitivities of inspectors. Nevertheless, the consistency across priors of the rankings of inspector sensitivities in the Mariano Melgar survey (Table 4) suggests that the relative sensitivities of the inspectors may be learned from the presence-absence data, even in the absence of reliable prior information.

Table 5: Simulation Mean (and SD) of Brier Scores Measuring the Effect of Factors on Infestation Estimation

Brier Scores for Infestation Estimates ($\ y_{\text{Sim}} - \hat{y}\ _2$)				
	Survey Assignment		Randomized Assignment	
Prior	Mean	(SD)	Mean	(SD)
B(1, 1)	21.2609	(0.2835)	21.0009	(0.2850)
B(5, 5)	22.3966	(0.2542)	22.1084	(0.2473)
B(6.5, 2)	20.9187	(0.2308)	20.7710	(0.3225)

Sample size: 50

The mean and standard deviation across simulations of the Brier scores, measuring how accurately we were able to estimate the infestation. Each cell (pair) represents a different combination of factors, where lower Brier scores are indicative of better estimates of the infestation under these factors.

Table 6: Simulation Mean (and SD) of Squared Norms Measuring the Effect of Factors on Accuracy of Sensitivity Estimation

Squared Norm of Inspector Sensitivity Estimates ($\ \beta_{\text{Sim}} - \hat{\beta}\ _2$)				
	Survey Assignment		Randomized Assignment	
Prior	Mean	(SD)	Mean	(SD)
B(1, 1)	1.4122	(0.0806)	1.2454	(0.1289)
B(5, 5)	1.8430	(0.0411)	1.7918	(0.0486)
B(6.5, 2)	0.6613	(0.0479)	0.6848	(0.0561)

Sample size: 50

The mean and standard deviation across simulations of the squared norm difference between our estimates of the inspector sensitivities and the generating inspector sensitivity parameters.

6 Discussion

We proposed a spatial model to analyze spatially-clustered presence-absence data that quantifies the amount of under-reporting and account for data missing not at random. The model allows us to capture the heterogeneity of the surveillance error across the individuals collecting the data. Applying our model to surveys for the presence of *T. infestans* in the district of Mariano Melgar in Arequipa, Peru, we identified two additional at-risk localities for treatment. Applying a simpler model, not accounting for the difference in participation between infested and non-infested households we found four additional at risk localities for treatment.

The willingness of infested households to participate affirm the Ministry’s local community outreach efforts and the willingness of communities severely affected by triatominae infestations to cooperate. We have identified these locality-level differences in participation previously in [27]. The Ministry of Health based its treatment decisions on estimates produced under the MAR assumption. Because we found a link between under-participation and lower rates of infestation, it is likely that these estimates were an overestimate of prevalence as evidenced by the estimates from our more complete MNAR model. Because of the strong entomological risk of re-infestation in triatominae insects, overestimates of prevalence may be of value to safeguard against the risk of re-infestation. However, because of the additional strain on public resources, these distinctions and assumptions should be clearly outlined to policymakers.

On simulated data, we showed that a hypothetical randomized assignment of inspectors only marginally improved the estimation of the infestation and inspector sensitivities. This similarity suggests that the assignment used to conduct the Mariano Melgar survey is not seriously susceptible to confounding issues in spite of some spatial correlation.

We found our model to produce similar interpolation of the risk of infestation in non-participating households than a more traditional regression model. Because we model the effect of surveillance error using an external Bernoulli random variable to the linear infestation risk, we believe the value of our model is that it is more easily interpreted. In our model, the true unobservable phenomenon of interest is modelled separately from the observed data. Policy decisions can then be made on this distinct quantity and false-negative probability is explicitly estimated.

Related work was done in the context of sociological surveys in [28], where researchers incorporated the covariates of the interviewers conducting the survey to help explain the variation in the data. These authors

found that questionnaires conducted regarding perception of social disorder in urban neighbourhoods were consistent across various interviewers. In this study, we found that posterior distributions of inspector detection abilities to depart strongly from the prior distribution. Although the overall posterior probability of detection was dependent on the particular prior specification, we found the rankings of inspector abilities to be consistent across a variety of prior distributions.

There are limitations to this study. First, despite accounting for the influence of the infestation on the participation and variations in participations between localities, we may be overlooking other covariates, which may also be key in mitigating the confounding. In our previous work, [24] we found the importance of covariates such as livestock and household building materials to be limited compared to the effect of the spatial correlation. In addition, our model for the spatial effect, Equation (3), is based heavily on thresholding squared distance and does not have an easily interpretable covariance [29].

Individual surveillance error models are useful for incorporating inspector-to-household labelling data in presence-absence estimates. The model proposed here not only quantifies the amount of under-reporting in survey data, but allows for the relative estimation of inspector quality. The infestation maps are produced efficiently at a fine resolution and account for non-participating households even missing non at random, which gives insight into the distribution of residual infestation post treatment.

6.1 Acknowledgements

The Chagas Disease Working Group in Arequipa includes Fernando Malaga Chavez, Karina Oppe Alvarez, Andy Catacora Rospigliossi, Claudia Patricia Mena Cornejo, Dr. Juan Cornejo del Carpio, Javier Quintanilla Calderon, Kate Levy, Malwina Niemierko, Victor Quispe Machaca, Jenny Ancca, Katty Borrini, Renzo Salazar Sanchez, Katty Borrini-Mayori, Danitza Pamo, Giovanna Moscoso, Lina Mollesaca, Maria Luz Hanco, and Renzo Salazar. We thank Dr. Christopher Paciorek and Dr. Tony Smith for insightful comments. We gratefully acknowledge the work of the following organizations that have organized and conducted the Chagas disease control program in Arequipa: Ministerio de Salud del Perú (MINSA), the Dirección General de Salud de las Personas (DGSP), the Estrategia Sanitaria Nacional de Prevención y Control de Enfermedades Metaxénicas y Otras Transmitidas por Vectores (ESNPCEMOTVS), the Dirección General de Salud Ambiental (DIGESA), the Gobierno Regional de Arequipa,

the Gerencia Regional de Salud de Arequipa (GRSA), the Pan American Health Organization (PAHO/OPS) and the Canadian International Development Agency (CIDA). We thank the district of Mariano Melgar for its participation in this study.

A Gibbs Sampler

We now outline the Gibbs sampler for implementing the model. Using the parameter expansion popularized in [30] for the probit link, the closed forms for all of the conditional distributions are known and in the form of common distributions. Modifications can be made for the logistic link by following [31]. The parameter expansion for the binary outcome y is implemented by introducing the continuous variable y_0 .

$$y_0 = u + t + \epsilon \quad (8)$$

$$y_{1,i} = \mathbb{1}_{\{y_{0,i} > 0\}} \quad (9)$$

where u is the conditionally auto-regressive model with precision $\mathbf{N}(0, k_u \mathbf{\Lambda})$, t is the intercept, and ϵ is the standard Gaussian. If the prior on t is given by $\mathbf{N}(\mu, \tau)$, where again τ is the precision, the prior on k_u is given by $\Gamma(k, \theta)$, where k and θ represent the scale and shape parameters, and the prior on each element of β is $\mathbf{B}(a, b)$. The conditional distributions for the model parameters are then given by the following:

1. $(k_{\mathbf{u}} | \mathbf{u}) \sim \Gamma\left(\frac{n-1}{2} + k, \frac{1}{2} \mathbf{u}^\top \mathbf{\Lambda} \mathbf{u}\right)$
2. $([\mathbf{u}, t]^\top | k_u, y_0) \sim \mathbf{N}\left(\begin{bmatrix} k_u \mathbf{\Lambda} + \mathbf{I} & 1 \\ 1^\top & n + \tau \end{bmatrix}^{-1} \begin{bmatrix} y_0 \\ 1^\top y_0 + \mu + \tau \end{bmatrix}, \begin{bmatrix} k_u \mathbf{\Lambda} + \mathbf{I} & 1 \\ 1^\top & n + \tau \end{bmatrix}\right)$
3. $(y_{0,i} | u_i, t, y_{1,i}) \sim \begin{cases} \mathbf{N}(u_i + t, 1 | y_{0,1} > 0) & \text{if } y_{1,i} = 1 \\ \mathbf{N}(u_i + t, 1 | y_{0,1} < 0) & \text{if } y_{1,i} = 0 \end{cases}$
4. $(y_{1,i} | u_i, \beta(i)) \sim \text{Bern}\left(p_i = \begin{cases} \frac{(1-\beta(i))\Phi(u_i+t)}{(1-\beta(i))\Phi(u_i+t)+(1-\Phi(u_i+t))} & \text{if } \mathbb{1}_{\text{NA}i} = 0 \\ \Phi(u_i + t) & \text{if } \mathbb{1}_{\text{NA}i} = 1 \end{cases}\right)$
5. $(\beta_i | y_{I_i}, z_{I_i}) \sim \mathbf{B}(\sum_{j \in I_i} y_j z_j + a, \sum_{j \in I_i} y_j (1 - z_j) + b)$

The Mariano Melgar analysis was performed in R on an Intel Core i7 processor clocked at 2.8 GHz, where one thousand iterations of the Markov chain were performed in 299.77 seconds. The size of the precision matrix in this work was 12,050-by-12,050 with 0.18% sparsity. We found convergence to be consistent irrespective of the starting point. We found the slowest mixing and most autocorrelated variable in the chain to be k_u , due to the strong dependence between u and k_u in the sampling scheme used. Based on the samples of k_u produced by the chain, we recommend discarding the first 10,000 samples as burn-in. After burn-in, we recommend thinning the samples and retaining only every tenth sample.

References

- [1] M.G. Teixeira, M.L. Barreto, M.C. Costa, L.D. Ferreira, P.F. Vasconcelos, and S. Cairncross. Dynamics of dengue virus circulation: a silent epidemic in a complex urban area. *Tropical Medicine & International Health*, 7(9):757–762, 2002.
- [2] E. Lars, B.J. Beaty, A.C. Morrison, and T.W. Scott. Proactive vector control strategies and improved monitoring and evaluation practices for Dengue prevention. *Journal of Medical Entomology*, 46(6):1245–1255, 2009.
- [3] L. Sanchez, J. Cortinas, O. Pelaez, H. Gutierrez, D. Concepcion, and P. Van der Stuyft. Breteau index threshold levels indicating risk for dengue transmission in areas with low aedes infestation. *Tropical Medicine & International Health*, 15(2):173–175, 2010.
- [4] J.F. Trape, E. Lefebvre-Zante, F. Legros, G. Ndiaye, H. Bouganali, P. Druilhe, and G. Salem. Vector density gradients and the epidemiology of urban malaria in Dakar, Senegal. *American Journal of Tropical Medicine and Hygiene*, 47(2):181–189, 1992.
- [5] S. Wang, C. Lengeler, T.A. Smith, P. Vounatsou, D.A. Diadie, X. Pritroipa, N. Convelbo, M. Kientga, and M. Tanner. Rapid urban malaria appraisal (RUMA) I: epidemiology of urban malaria in Ouagadougou. *Malaria Journal*, 5(43), 2005.
- [6] B. Matthys, G. Vounatsou, A.B. Tschannen, E.G. Becket, L. Gosoni, G. Cissé, M. Tanner, E.K. N’goran, and J. Utzinger. Urban farming and malaria risk factors in a medium-sized town in Côte d’Ivoire. *American Journal of Tropical Medicine and Hygiene*, 75(6):1223–1231, 2006.

- [7] H.J. Corrasco, A. Torrellas, C. García, M. Segovia, and M.D. Felician-geli. Risk of *Trypanosoma cruzi* I (Kinetoplastida: Trypanosomatidae) transmission by *Panstrongylus geniculatus* (Hemiptera: Reduviidae) in Caracas (Metropolitan District) and neighboring States, Venezuela. *International Journal for Parasitology*, 35(13):1379–1384, 2005.
- [8] J.C. Dias, A.C. Silveira, and C.J. Schofield. The impact of Chagas disease control in Latin America: a review. *Memórias do Instituto Oswaldo Cruz*, 97(5):603–612, 2002.
- [9] F. Guhl. Chagas disease in andean countries. *Mem. Inst. Oswaldo Cruz*, 1:29–38, 2007.
- [10] M.Z. Levy, N.M. Bowman, V. Kawai, L.A. Waller, J.G. Cornejo del Carpio, E. Cordova Benzaquen, R.H. Gilman, and C. Bern. Periurban *Trypanosoma cruzi*-infected *Triatoma infestans*, Arequipa, Peru. *Emerging Infectious Diseases*, 12(9):1345–1352, 2006.
- [11] MD Germano, G. R. Acevedo, G. A.M. Cueto, AC Toloza, CV Vassena, and MI Picollo. New findings of insecticide resistance in *Triatoma infestans* (Heteroptera: Reduviidae) from the Gran Chaco. *Journal of medical entomology*, 47(6):1077–1081, 2010.
- [12] Frdric Lardeux, Stphanie Depickre, Stphane Duchon, and Tamara Chavez. Insecticide resistance of *triatoma infestans* (hemiptera, reduviidae) vector of chagas disease in bolivia. *Trop Med Int Health*, 15(9):1037–1048, Sep 2010.
- [13] M.Z. Levy, F.S. Malaga Chavez, J.G. Cornejo Del Carpio, D.A. Vilhena, F.E. McKenzie, and J.B. Plotkin. Rational spatio-temporal strategies for controlling a Chagas disease vector in urban environments. *Journal of the Royal Society Interface*, 7:1061–1070, 2010.
- [14] C. Barbu, E. Dumonteil, and S. Gourbière. Evaluation of spatially target strategies to control non-domiciliated *triatoma dimidiata* vector of Chagas disease. *PLoS Neglected Tropical Diseases*, 5:e1045, 2011.
- [15] Roderick J Little, Ralph D’Agostino, Michael L Cohen, Kay Dickersin, Scott S Emerson, John T Farrar, Constantine Frangakis, Joseph W Hogan, Geert Molenberghs, and Susan A Murphy. The prevention and treatment of missing data in clinical trials. *New England Journal of Medicine*, 367(14):1355–1360, 2012.

- [16] Nhu D Le, Weimin Sun, and James V Zidek. Bayesian multivariate spatial interpolation with data missing by design. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 59(2):501–510, 1997.
- [17] Christel Faes, JT Ormerod, and MP Wand. Variational bayesian inference for parametric and nonparametric regression with missing data. *Journal of the American Statistical Association*, 106(495), 2011.
- [18] Emily L Kang and Noel Cressie. Bayesian inference for the spatial random effects model. *Journal of the American Statistical Association*, 106(495):972–983, 2011.
- [19] Google. Google earth. Accessed on 17/08/2009.
- [20] S. Banerjee, M. M. Wall, and B. P. Carlin. Frailty modeling for spatially correlated survival data, with application to infant mortality in minnesota. *Biostatistics*, 4:123–142, 2003.
- [21] S. Banerjee, B. P. Carlin, , and A. E. Gelfand. *Hierarchical Modeling and Analysis for Spatial Data*. Chapman & Hall, first edition, 2004.
- [22] J. Besag, J. York, and A. Mollié. Bayesian image restoration with two applications in spatial statistics. *Ann. Inst. Statist. Math.*, 43(1):1–59, 1991.
- [23] H. Rue and L. Held. *Gaussian Markov Random Fields Theory and Applications*. Chapman & Hall, first edition, 2005.
- [24] C. Barbu, A. Hong, J. M. Manne, D. Small, J. E. Calderón, K. Sethuraman, V. Quispe-Machaca, J. Ancca-Juárez, J. G. Cornejo del Carpio, F. S. Chavez, et al. The effects of city streets on an urban disease vector. *PLoS Computational Biology*, page e1002801, 2013.
- [25] C. Paciorek. Computational techniques for spatial logistic regression with large data sets. *Computational Statistics & Data Analysis*, 51(8):3631–3653, 2007.
- [26] G.W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78:1–3, 1950.
- [27] Alison M Buttenheim, Valerie Paz-Soldan, Corentin Barbu, Christine Skovira, Javier Quintanilla Calderón, Lina Margot Mollesaca Riveros, Juan Oswaldo Cornejo, Dylan S Small, Christina Bicchieri, Cesar

- Naquira, et al. Is participation contagious? evidence from a household vector control campaign in urban peru. *Journal of Epidemiology and Community Health*, 68:103–109, 2014.
- [28] F. Casas-Cordero, F. Kreuter, Y. Wang, and S. Babey. Assessing the measurement error properties of interview observations of neighborhood characteristics. *J. R. Statist. Soc. A*, 176(1):227–249, 2013.
- [29] D. Pickard. Asymptotic inference for an Ising lattice. ii. *Adv. Appl. Prob.*, 9:476–501, 1977.
- [30] J. H. Albert and S. Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–679, 1993.
- [31] C.C. Holmes and L. Held. Bayesian auxiliary variable models for binary and polychotomous regression. *Bayesian Analysis*, 1(1):145–168, 2006.