

Inference of malaria transmission dynamics under varying assumptions about the spatial scale of transmission: a modelling study in three elimination settings

Isobel Routledge (✉ isobel.routledge@ucsf.edu)

University of California, San Francisco

Samir Bhatt

Imperial College London

Research Article

Keywords: Individual-level, parameterisations, Euclidian distance, genetic distance, accessibility matrices

Posted Date: November 25th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-107010/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Scientific Reports on July 14th, 2021. See the published version at <https://doi.org/10.1038/s41598-021-93238-0>.

Inference of malaria transmission dynamics under varying assumptions about the spatial scale of transmission: a modelling study in three elimination settings

Isobel Routledge^{1*}, Samir Bhatt²

1. University of California, San Francisco 2. Imperial College London

*Correspondence to: isobel.routledge@ucsf.edu

Abstract

Background

Individual-level geographic information about malaria cases, such as the GPS coordinates of residence or health facility, is often collected as part of surveillance in near-elimination settings, but could be more effectively utilised to infer transmission dynamics, in conjunction with additional information such as symptom onset time and genetic distance. However, in the absence of data about the flow of parasites between populations, the spatial scale of malaria transmission is often not clear. As a result, it is important to understand the impact of varying assumptions about the spatial scale of transmission on key metrics of malaria transmission, such as reproduction numbers.

Methods

We developed a method which allows the flexible integration of distance metrics (such as Euclidian distance, genetic distance or accessibility matrices) with temporal information into a single inference framework to infer malaria reproduction numbers. Twelve scenarios were defined, representing different assumptions about the likelihood of transmission occurring over different geographic distances and likelihood of missing infections (as well as high and low amounts of uncertainty in this estimate). These scenarios were applied to four individual level datasets from malaria eliminating contexts to estimate individual reproduction numbers and how they varied over space and time.

Results

Model comparison suggested that including spatial information improved models as measured by $\Delta AICc$, compared to time only results. Across scenarios and across datasets, including spatial information tended to increase the seasonality of temporal patterns in reproduction numbers and reduced noise in the temporal distribution of reproduction numbers. The best performing parameterisations assumed long-range transmission (>200km) was possible.

Conclusions

Our approach is flexible and provides the potential to incorporate other sources of information which can be converted into distance or adjacency matrices such as travel times or molecular markers.

35 Introduction

36 Individual-level disease surveillance data, collected routinely and as part of outbreak response, capture
37 a wealth of information which could improve measurements of transmission and its spatiotemporal
38 variation, in turn informing the design of epidemiological interventions¹. Geo-located health facility or
39 residence data are increasingly collected as part of surveillance for diseases such as malaria², but could
40 be more effectively utilised to infer transmission dynamics, in conjunction with additional information
41 such as symptom onset time and genetic distance³. However, challenges exist in making use of these
42 diverse data sources and leveraging the information they contain within a single inference framework.
43 This is particularly true of endemic diseases such as malaria, where individual level data are increasingly
44 collected in moderate-low transmission or elimination settings.

45 Malaria transmission is shaped by processes occurring on a wide range of spatial scales. In the absence
46 of human mobility, transmission is limited to the range of the mosquito vector, however human
47 movement, ranging from regular commutes to rare large scale migration events, can import parasites
48 into new areas provided competent vectors are present⁴⁻⁷. As a result, the spatial location of individual
49 cases can provide useful information in inferring transmission dynamics when combined with additional
50 forms of information, such as temporal and molecular data. Furthermore, in elimination settings malaria
51 transmission is thought to take on epidemic dynamics⁸, meaning the importance of space and highly
52 dynamic factors such as human movement patterns becomes more relevant. However, in the absence
53 of data about the flow of parasites between populations, the spatial scale of malaria transmission is
54 often not clear. As a result, it is important to understand the impact of varying assumptions about the
55 spatial scale of transmission on key metrics of malaria transmission, such as reproduction numbers. In
56 many contexts, not all cases may be observed within a surveillance system. Missing cases further
57 complicate inference, as without information about how likely cases are to be missing, ambiguity can

58 exist as to whether long-range transmission occurred or whether cases were infected by a closer,
59 unobserved source of infection.

60 To explore the impact of including distance measure and assumptions about their relationship to
61 transmission likelihood, we developed a flexible framework to incorporate pairwise distances (for
62 example Euclidian distances, travel times, or any quantifiable distance matrix) into our previously
63 published inference framework⁹ to estimate individual reproduction numbers and explore the impact of
64 varying assumptions about missing cases and the spatial kernel on results, as well as determining the
65 feasibility of inferring the distance kernel and amount of missing cases from surveillance data. We
66 defined twelve scenarios representing different assumptions about the likelihood of transmission
67 occurring over different geographic distances and likelihood of missing infections (as well as high and
68 low amounts of uncertainty in this estimate). These scenarios were applied to four individual level
69 datasets from malaria eliminating contexts to estimate individual reproduction numbers and how they
70 varied over space and time. We used two simple spatial kernels describing the relationship between
71 Euclidian distance between residences and likelihood of transmission occurring, to explore various
72 assumptions about the relationship between locations of cases and likelihood of transmission occurring
73 between them, as well as the impact of unobserved cases. We find the best performing models by
74 second order AIC ($\Delta AICc$)¹⁰ have weakly informative priors on the likelihood of unobserved sources of
75 infection and assume . However, we find there can be issues of parameter identifiability, which become
76 increasingly relevant when there are not enough data available about key parameters in the model.

77 Results

78 We developed a framework to integrate distance information into a previously published inference
79 framework^{9,11} which uses the time of symptom onset to infer reproduction numbers and their
80 spatiotemporal variation. We then tested the impact of varying assumptions about the relationship

81 between location of cases and the likelihood of transmission as well as the impact of unobserved
 82 infection as modelled by competing edges, ϵ , considering twelve scenarios (**Table 1**), and applying them
 83 to four line-list datasets from China (*P. vivax* and *P. falciparum*, analysed separately), El Salvador (*P.*
 84 *vivax*) and Eswatini (*P. falciparum*) using Exponential and Gaussian spatial kernels, described in the
 85 methods section of this paper.

86 *Table 1: Table illustrating the different scenarios and corresponding parameter values tested in scenario*
 87 *analysis*

Scenario Description	Scenario	Beta (fixed)	Epsilon (prior)
Human movement unlikely, most movement under 10km Missing cases more likely (but very uncertain)	1	Gaussian = 0.005 Exponential =0.1	Mean = 0.1 SD = 1
Human movement unlikely, most movement under 10km Missing cases more likely (confident)	2	Gaussian = 0.005 Exponential =0.1	Mean = 0.1 SD = 0.001
Human movement unlikely, most movement under 10km Missing cases less likely (but very uncertain)	3	Gaussian = 0.005 Exponential =0.1	Mean = 0.001 SD = 1
Human movement unlikely, most movement under 10km Missing cases less likely (confident)	4	Gaussian = 0.005 Exponential =0.1	Mean =0.001 SD =0.001
Moderate human movement, most movement under 50km Missing cases more likely (but very uncertain)	5	Gaussian = 0.001 Exponential =0.02	Mean = 0.1 SD = 1
Moderate human movement, most movement under 50km Missing cases more likely (confident)	6	Gaussian = 0.001 Exponential =0.02	Mean = 0.1 SD = 0.001
Moderate human movement, most movement under 50km Missing cases less likely (but very uncertain)	7	Gaussian = 0.001 Exponential =0.02	Mean = 0.001 SD = 1
Moderate human movement, most movement under 50km Missing cases less likely (confident)	8	Gaussian = 0.001 Exponential =0.02	Mean =0.001 SD =0.001
Longer range human movement likely Missing cases more likely (but very uncertain)	9	Gaussian = 0.0001 Exponential =0.01	Mean = 0.1 SD = 1
Longer range human movement likely Missing cases more likely (confident)	10	Gaussian = 0.0001 Exponential =0.01	Mean = 0.1 SD = 0.001
Longer range human movement likely Missing cases less likely (but very uncertain)	11	Gaussian = 0.0001 Exponential =0.01	Mean = 0.001 SD = 1
Longer range human movement likely Missing cases less likely (certain)	12	Gaussian = 0.0001 Exponential =0.01	Mean =0.001 SD =0.001

88

89 **Results of model comparison by $\Delta AICc$ across different scenarios**

90 When $\Delta AICc$ scores were used to compare model results, all models which included distance had lower
 91 (and therefore better) $\Delta AICc$ scores than models which only included only time (**Table 2 and**
 92 **Supplementary Table 1**). In addition, exponential kernels consistently outperformed equivalent
 93 scenarios using Gaussian kernels (**Supplementary Table 1**). Two scenarios consistently performed best
 94 as measured by $\Delta AICc$, namely Scenario 9 (El Salvador and Swaziland) and Scenario 11 (China, *P. vivax*
 95 and *P. falciparum*). Both scenarios assume longer range human movement likely and impose a smaller
 96 penalty on cases occurring larger distances. These scenarios also allow variation in epsilon edge values
 97 and use a very weakly informative prior on Epsilon edges, but with a different mean (0.1 for Scenario 9,
 98 0.001 for Scenario 11). These results also return smaller mean R_c results than time-only versions of the
 99 model (**Figures 1– 4**)

100 *Table 2: Summary of $\Delta AICc$ results*

Dataset	Best Model(s), by $\Delta AICc$	Akaike Weight
Swaziland (Eswatini)	Scenario 9, Exponential	1
El Salvador	Scenario 9, Exponential Scenario 11, Exponential	0.621540909785805 0.37845909
China <i>P. vivax</i>	Scenario 11, Exponential	1
China <i>P. falciparum</i>	Scenario 11, Exponential	1

101

102 **R_c estimates under different scenarios**

103 Across all datasets, large differences in R_c estimates were found depending on both ϵ and
 104 β parameters. When β is higher, the assumption is that there is little movement of parasites within the
 105 country and therefore cases with residential addresses which are far away are unlikely to have infected
 106 each other. When this is the case and we assume there are unobserved sources of infection (either

107 through a strongly informative prior on ε with mean 0.1, or an uninformative prior with a lower mean),
108 then R_c values are very low. However if we assume there are little or no unobserved sources of
109 infection, but continue to make restrictive assumptions about space, then most R_c very low but in the
110 localities where there are cases we estimate much higher R_c values as there are no other possible
111 infectors within a reasonable time and/or spatial area. This is illustrated in Figures 1 - 4.

112 **[Insert Figures 1-4 approximately here]**

113
114 When looking at the spatial patterns of R_c estimates under different scenarios several trends are seen
115 across all datasets (**Figures 5-8**). Scenario 4 is particularly interesting to note because this scenario
116 considers the most restrictive assumptions, both about space and unobserved sources of infection.
117 Across datasets, Scenario 4 results in increased focality and higher R_c s within these foci, but in
118 comparison lower R_c s in other areas. All of the best scenarios as measured by $\Delta AICc$ resulted in small R_c
119 estimates, but where comparably larger R_c estimates were estimated, they were in localities identified
120 as foci.

121
122 **[Insert Figures 5-8 approximately here]**

123 For the line-list dataset from El Salvador, within the range of values explored in the sensitivity
124 analysis (**Table 3**), regardless of how informative the prior was for either β , the distance shaping
125 function, or for ϵ , the epsilon edge, β was always estimated as whatever the mean of the prior was
126 set as between the prior mean values of $1e-4$ and $1e-2$ (**Figure 9**). However, when the mean value
127 was set at 0.1, the estimated parameter converged at a slightly lower value of 0.075, with the
128 exception of when the prior for ϵ was very low (all priors with mean ϵ of $1e-10$ and also the more
129 informative priors with mean $1e-5$, when standard deviation was $1e-4$). R_c is strongly shaped by the
130 value of ϵ , with higher values of ϵ returning lower values of R_c , however R_c also declined with
131 increasing values of β .

132 Very similar patterns to El Salvador were observed in the sensitivity analysis of the Eswatini dataset.
133 Again, regardless of how informative the prior was for either ϵ or β , β was always estimated as
134 whatever the mean of the prior was set as between the prior mean values of $1e-4$ and $1e-2$ (**Figure**
135 **10**). However, when the mean value was set at 0.1, the estimated parameter converged at a slightly
136 lower value of 0.075, with the exception of when the prior for ϵ was very low (all priors with mean ϵ
137 of $1e-10$ and also the more informative priors with mean $1e-5$, when standard deviation was $1e-4$).
138 Unlike El Salvador, for Eswatini, at higher values of ϵ (0.5 and 0.1) there are stark declines in R_c with
139 increasing β .

140 For both *P. vivax* and *P. falciparum* datasets from China, within the parameter range explored in the
141 sensitivity analysis, regardless of how informative the prior was for either β , the distance shaping
142 function, or ϵ , the epsilon edge, β was always estimated as whatever the mean of the prior was set
143 as (**Figures 11 and 12**), suggesting a lack of identifiability or information within the data. When
144 estimating R_c , and interesting interacting effect of ϵ (missing or unobserved infections) and β
145 (distance) was seen. When β is low, although lower values of ϵ produce slightly higher
146 mean R_c values, the difference in R_c estimates with varying prior values for ϵ is much smaller than
147 when β is a higher value. In other words, when the prior for ϵ is low, $1e-10$, R_c estimates do not vary

148 as β changes, however when the prior for ϵ is much higher, then increasing β from $1e-4$ to 0.1
 149 reduces R_c estimates (from 0.21 to 0.01 for *P. vivax*).

150 *Table 3: Different parameters considered in sensitivity analysis. Note all combinations of each parameter were*
 151 *considered.*

ϵ mean	ϵ SD	β mean (Gaussian)	β mean (Exponential)	β SD
1e-10	0.0001	0.00001	0.0001	0.0001
1e-5	0.001	0.0001	0.001	0.001
1e-3	0.01	0.001	0.01	0.01
1e-2	0.05	0.01	0.1	0.05
1e-1	0.1			0.1
0.5				

152

153 **[Insert Figures 9-12 approximately here]**

154

155 Discussion

156 We developed an approach to estimate malaria transmission network properties, which allows the
157 flexible integration of distance metrics, such as Euclidian distances or travel times, with temporal
158 information within a single inference framework. Twelve scenarios and corresponding parameter
159 values were defined which represented a) varying likelihoods of transmission over different
160 distances and b) varying likelihoods of missing infections (as well as high and low confidence in this
161 estimate). These scenarios were applied to four individual level datasets from malaria eliminating
162 contexts and using two different spatial kernels. The estimated R_c values, their spatial and temporal
163 distribution and the $\Delta AICc$ /Akaike weights for each model were compared alongside a time only
164 model. These results suggest that including spatial information improved models as measured by
165 AIC, compared to time only results. The prior values for both the distance function and epsilon value
166 have very strong impacts on the estimated R_c , although relative temporal trends tend to stay
167 consistent.

168 For all datasets considered, all model versions which used geographic information had lower $\Delta AICc$
169 values than the time only model. Based on the Akaike Weights and $\Delta AICc$ values for each model,
170 large differences in $\Delta AICc$ were seen between different scenarios. Scenarios 9 and 11 produced the
171 lowest $\Delta AICc$ values. These were parameterisations which penalised long range transmission the
172 least where and the prior on epsilon edges was only weakly informative. These parameterisations
173 also return much lower reproduction numbers than using time alone.

174 Exponential Kernels consistently outperformed Gaussian kernels as measured by $\Delta AICc$. Although
175 classic models of dispersion are as a diffusion process with Gaussian displacement, more leptokurtic
176 or “fatter-tailed” probability distributions, where more of the probability density is concentrated in
177 the tails of the function, are often found to better represent empirical dispersal patterns than
178 traditional Gaussian kernels¹². This “fatter-tail” in the exponential can be seen in Figures 13 - 15.

179 However, there are many limitations to using ΔAICc in model comparison, particularly when
180 estimation of some of the parameters are being carried out within a Bayesian context. We do not fix
181 α_{ij} nor do we fix epsilon, but we do define priors and maximise the posterior rather than the log
182 likelihood. Therefore, we are comparing negative log likelihoods from a maximised posterior,
183 meaning we are not considering the information included in the prior. In addition, many α_{ij} values
184 shrink to zero, however are still counted as parameters in the AIC estimation. Therefore, there is no
185 recognition of which versions of the model produce fewer non-zero parameters. Whilst this
186 difference in AIC is interesting to note, I would argue the broader trends in how R_c varies over time
187 and space with different assumptions about both the spatial kernel and the number of unobserved
188 sources of infection are more important to consider.

189 An interesting pattern which was noted across scenarios and across datasets was how including
190 spatial information in the likelihood tended to increase the seasonality of temporal patterns in
191 reproduction numbers and reduced noise in the temporal distribution of reproduction numbers. This
192 could be suggestive of importation events leading to localised infections. Scenario 4 is also an
193 interesting set of assumptions to consider as it assumes cases generally only infect cases near them
194 and that unobserved cases of infection are unlikely. Under this assumption foci of infection are very
195 clear and clear “sources” of infection.

196 The results of the sensitivity analysis reveal interesting differences between the different datasets
197 and contexts contained in this dataset. For both El Salvador and Eswatini, which are both small
198 countries (El Salvador has an area of 21,041 km² and Eswatini 17,364 km²), at higher mean priors for
199 β , the model converged on an estimate for β which was informed by the data. This was not the case
200 for the dataset from China, which represents a much larger area geographically and where dynamics
201 are likely to be strongly driven by importation. Given that for the kernels we used in this analysis,
202 increasing values of β lead to more restrictive assumptions about the scale of transmission, perhaps
203 this difference is due to the different spatial scales at which the analysis was being carried out.

204 There are several limitations to this approach and analysis. Firstly, there is a potential lack of
205 identifiability between ϵ , the epsilon edge, and β , the shaping parameter of the spatial kernel. To
206 give an intuitive example, say two cases occurred 50km from each other in space within a
207 reasonable timeframe of symptom onset times for transmission to have occurred. Without strong
208 prior information about what the spatial kernel may be, and/or how likely cases are to have an
209 external source of infection, it is not clear whether these cases are linked by transmission (and there
210 is some human travel/parasite movement, modelled by a less restrictive spatial kernel) or whether
211 there are unobserved source(s) of infection leading to both cases. This is also exemplified in the
212 results of the sensitivity analysis, where the mean of the prior for beta strongly shapes the final
213 estimate of beta, and the epsilon value also shapes beta.

214 In the absence of reliable information about either of these values, strong assumptions must be
215 made about either/both the likelihood of cases being infected by unobserved sources of infection
216 and the relationship between distance and. Similar approaches¹³ recommend fixing the kernel
217 shaping parameter, and indeed approaches from others have also noted problems with
218 unconstrained distance kernels in space-time diffusion modelling (Swapnil Mishra, personal
219 correspondence). One potential way to address this is divide epsilon edge by the distance parameter
220 $\frac{\epsilon}{\beta}$, thereby linking the two parameters and thereby penalising increases in β .

221 Indeed, for similar approaches analysing the diffusion of twitter hashtags, it was recommended to fix
222 the parameter beta, and the authors acknowledged potential challenges in estimating this
223 parameter. Whilst the temporal aspect is not fixed, I view the utility in this method in excluding or
224 penalising improbable transmission links between far away cases, rather than as a way of trying to
225 determine what the spatial relationship between cases is for malaria transmission, or determining
226 the relative contribution of space to malaria transmission.

227 An additional approach which could alleviate this problem is to collect internal travel history as part
228 of surveillance in future data collection efforts. This may help tease apart the relationship between

229 space and transmission. There also may be regions where there is more information to parameterise
230 both the spatial scales of transmission and the likelihood of cases being unobserved (for example
231 through looking at reporting rates, rates of relapse in the case of *P. vivax*, and prevalence of
232 asymptomatic infection).

233 Secondly, our approach was designed for application to near elimination and elimination settings,
234 where surveillance and case management is very strong, numbers of cases are small, and therefore
235 there is less overlap in potential infector/infectees, and changes in transmission are more apparent.
236 If applying these approaches to contexts which are less far along the journey to elimination, the
237 issue of identifiability may be even more of an issue as one cannot reasonably assume/fix epsilon
238 edges to be a very small number. Asymptomatic infection will likely be more important to consider,
239 more sophisticated methods to deal with missing cases will be required. There also will likely be a
240 weaker signal in space and time, which may require the integration of additional information such as
241 genetic distance. There also will be a transmission level above which these methods will no longer
242 be useful, although we do not know what this exact level is.

243 Finally, due to there being no “ground truth” it is hard to rigorously compare model performance.
244 $\Delta AICc$ and Akaike Weights are standard, however as mentioned previously, there are important
245 limitations in using these metrics for model comparison. A useful future step would be to analyse
246 simulated line-lists which are spatially explicit to investigate the impact of varying parameter values
247 and the interaction between the shaping parameter of the spatial kernel and epsilon. Spatially
248 explicit simulations may also reveal how tolerant the method is to missingness.

249 Currently, missing cases are dealt with in a relatively simple way, under the assumption that in the
250 elimination settings used here, surveillance and control have been strong for an extended period of
251 time as to ensure small case numbers and low prevalence of asymptomatic parasitaemia, and that
252 the contribution of missing cases is small enough to be represented as a competing hazard.
253 However, if missingness was biased, it is not clear how strongly this would affect results. Further

254 simulations which model different forms of missing data/sampling schemes would be useful to
255 reveal the potential impact of non-random missing data. These simulations could also model
256 different sources of unobserved infection – for example missing cases caused by relapse of dormant
257 *P. vivax*, unreported cases or asymptomatic infection.

258 Many methods used to model and represent space and mobility have not been tested here due to
259 the issues of identifiability seen even in simple models of space. Gravity, radiation¹⁴, and friction
260 surfaces¹⁵ are all potentially useful models of how space may affect the likelihood of transmission. As
261 mosquitoes have a limited range and lifespan, developing better data and models of human
262 movement, and how it varies in different cultural contexts and between different demographic
263 groups, will provide useful information to appropriately parameterise and design the spatial
264 component of the model.

265 Although the prior for the shaping parameter of the serial interval was selected under the
266 assumption that the majority of cases are treated in a timely manner, In this analysis we have not
267 explicitly utilised information about the time and location of treatment, although this is available in
268 some contexts. This may be useful information to constrain the potential time window of infection
269 occurring, as detailed information about infectivity and gametocyte carriage following treatment
270 with anti-malarials is available¹⁴, although sub-optimal dosage, compliance and resistance have been
271 associated with differing outcomes and therefore having additional information about treatment
272 and prevalence of resistance would also be useful.

273 Another avenue for future work would be to adapt the approach to incorporate further sources of
274 information, such as genetic markers of similarity between parasites. For our approach to be useful
275 in contexts which are not at or within a few years of elimination, incorporation of additional
276 information into the inference framework will be required. This could be carried out either directly
277 by incorporating an additional term or function in the likelihood or indirectly through informing the
278 value of parameters and allowing them to vary between individuals. Previous work within the

279 machine learning and network analysis community has successfully integrated diverse sources of
280 information about texts such as language and similarity of context into very similar algorithms to the
281 one presented here¹³.

282 Conclusion

283 Increasingly, line-list data contain spatial and other forms of information. Developing rigorous
284 approaches to leverage the information contained within these diverse datasets will increasingly be
285 useful in malaria surveillance and epidemiology^{3,5,17} and developing a framework which flexibly
286 takes on different forms of data within an integrated inference framework is a key aspect of this.
287 There may be more useful information contained in genetic, and or travel, mobility data. However,
288 as we have seen there can be issues of identifiability, which becomes increasingly relevant when
289 there is not enough data available about key parameters in the model. Finding ways for leveraging
290 multiple datasets, understanding their relationships, how they can enhance info contained in others,
291 or used to build consensus is important.

292 We developed and tested an algorithm which flexibly allows the incorporation of distance or
293 adjacency matrices describing the distance or connectivity between cases. This was applied to
294 individual malaria case data from four eliminating and very low transmission contexts and a detailed
295 sensitivity analysis was carried out. The results of these analyses suggest that including space
296 improves model performance as measured by $\Delta AICc$, and that, for the contexts considered here, the
297 best performing models produce lower reproduction estimates than using temporal information
298 only, likely in part due to estimating more unobserved sources of infection. However, this
299 conclusion would be strengthened by more in-depth simulation studies. The approach presented
300 here could be adapted to many different datasets and contexts, however issues of identifiability
301 must be considered. The utility of this approach would be strengthened with further development of
302 the methods of modelling unobserved sources of infection. Our results also make it clear that in

303 many contexts that additional information sources may be required such as genetic or serological
304 data.

305 Methods

306 Data

307 *The Kingdom of Eswatini*

308 This dataset, previously analysed by Reiner and colleagues¹⁸ captures malaria cases recorded by the
309 National Malaria Elimination Programme in the Kingdom of Eswatini (formally known as Swaziland)
310 between January 2010 and June 2014. For each case detected during this time (N= 1373), case
311 investigation was carried out. For each case the following were collected: GPS coordinates of
312 household location, demographic information (age, occupation and sex), use of malaria prevention
313 interventions such as long-lasting insecticide treated bednets (LLINs), and date of symptom onset,
314 diagnosis and treatment, as well as travel history. Based on travel history cases were defined as
315 locally acquired, imported. For a small number of cases (N=58) the local/imported status was
316 determined “unknown”. For the purposes of this analysis, these cases were treated the same as local
317 cases, i.e. they were assumed to have potentially been infected by other cases in the dataset and/or
318 been infectors themselves.

319 *China*

320 This dataset consists of individual-level case data for all confirmed and probable cases reported in
321 China between 2011 and 2016^{9,19} (Table 4 and Table 5). The data consist of an individual identifier,
322 date of symptom onset, date of diagnosis and date of treatment, as well as the geolocated address
323 of residence and health facility. If the suspected location of infection was in China and not in the
324 same district, then the presumed location of infection was also included in the dataset. Demographic
325 information such as age and sex were also collected. For the analysis, data were separated into *P*.

326 *falciparum* and *P. vivax* cases. *P. malariae* (N=252) and *P. ovale* (N=822) were reported but excluded
 327 from the analysis due to the lower public health concern of these species.

328

329 *Table 4: Cases by diagnosis type (probable and confirmed) and species across China*

	Mixed infection	<i>P. falciparum</i>	<i>P. malariae</i>	<i>P. ovale</i>	<i>P. vivax</i>	Untyped
Confirmed	260	11830	252	822	6631	87
Probable	0	176	0	0	693	311

330

331 *Table 5: Cases by imported/local status and species across China*

	Mixed infection	<i>P. falciparum</i>	<i>P. malariae</i>	<i>P. ovale</i>	<i>P. vivax</i>	Untyped
Local	5	92	4	1	1711	95
Imported	255	11914	248	821	5613	303

332

333 *El Salvador*

334 This dataset consists of all confirmed cases of malaria recorded by the Ministry of Health in El
 335 Salvador between 2010 and the first two months of 2016¹¹ (N= 91 cases, of which 30 imported, 6 *P.*
 336 *falciparum*, 85 *P. vivax*). For each case, the date of symptom onset was recorded. Residential address
 337 was available for all but two cases. For these cases, the location was available at the *municipio*, or
 338 municipality level, and the coordinates of the centroid of the municipality (which for both were
 339 cities) were used as the geo-location. Two cases had addresses listed outside of El Salvador, both of
 340 which were located in Guatemala. All cases within El Salvador with full addresses (N=85) were
 341 georeferenced by latitude and longitude to *caserío/lotificación* level, which is approximately
 342 neighbourhood or hamlet level.

343 *Transmission model specifics*

344 In order to incorporate pairwise distance metrics, we extended our previously published algorithm
 345 applied to Yunnan Province, China⁹ by introducing a second function, f_2 , which describes the
 346 relationship between space (or distance of any kind) and likelihood of transmission. An appropriate
 347 function such as a Gaussian kernel is defined and the parameter(s) shaping that distribution, β , are

348 either fixed, or given a prior distribution and estimated from the data. Multiplied, together, this
 349 returns a single function:

350 $f(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = f_1(t_i | t_j; \alpha_{i,j}) \times f_2(x_i | x_j; \beta)$ (1)

351 Determined by times t , spatial locations x , transmission rates α , spatial parameter(s) β .

352 As before, the hazard is defined as the pairwise likelihood divided by the survival term:

353 $H = \frac{f(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta)}{S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta)}$ (2)

354 To derive the survival function, one integrates across all distances and times as follows:

355 $S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = (\int_0^\infty \int_0^{t_i-t_j} f_1(t_i | t_j; \alpha_{i,j}) f_2(x_i | x_j; \beta) dt dx$ (3)

356 The specific functions used in $f_1(t_i | t_j; \alpha_{i,j})$ and $f_2(x_i | x_j; \beta)$ will have large impacts on the
 357 outcomes of results and therefore the assumptions inherent in these choices must be made explicit
 358 and linked to the mechanisms of transmission.

359 To illustrate this approach by applying to several malaria line-lists, we used a shifted Rayleigh
 360 distribution to model serial interval distributions, $f_1(t_i | t_j; \alpha_{i,j})$. For the second part of the likelihood
 361 which model the relationship between space and the likelihood of transmission $f_2(x_i | x_j; \beta)$,
 362 Gaussian and Exponential diffusion kernels were used (Table 6).

363 *Table 6: Equations for f_1 , f_2 , hazard and survival for time -only, Gaussian and Exponential spatial kernels*

	$f_1(t_i t_j; \alpha_{i,j})$	$f_2(x_i x_j; \beta)$	Hazard	Survival
Exponential	$\alpha(t_i - t_j - \gamma)e^{-\frac{1}{2}\alpha(t_i-t_j-\gamma)}$	$e^{-\beta(x_i-x_j)}$	$\beta\alpha(t_i - t_j - \gamma)e^{-\beta(x_i-x_j)}$	$e^{-\frac{1}{2}\alpha(t_i-t_j-\gamma)} \frac{1}{\beta}$

Gaussian	$\alpha(t_i - t_j - \gamma)e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)}$	$e^{-\beta(x_i - x_j)^2}$	$\frac{2\sqrt{\beta}\alpha(t_i - t_j - \gamma)e^{-\beta(x_i - x_j)^2}}{\sqrt{\pi}}$	$e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)} \frac{\sqrt{\pi}}{2\sqrt{\beta}}$
-----------------	--	---------------------------	---	---

Time only	$\alpha(t_i - t_j - \gamma)e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)}$	n/a	$\alpha(t_i - t_j - \gamma)$	$e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)}$
------------------	--	-----	------------------------------	--

364

365

366 Using a shifted Rayleigh distribution and an exponential kernel the pairwise likelihood of a case

367 showing symptoms at t_i and at residence location x_i being infected by a case showing symptoms at

368 time t_j and at residence location x_j , becomes

$$f(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = \alpha(t_i - t_j - \gamma)e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)} e^{-\beta(x_i - x_j)^2} \quad (4)$$

369 The survival term simplifies to:

$$S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)} \frac{1}{\beta} \quad (5)$$

371 And the hazard simplifies to:

$$H(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = \beta\alpha(t_i - t_j - \gamma)e^{-\beta(x_i - x_j)^2} \quad (6)$$

373 For the Gaussian function, the pairwise likelihood of a case showing symptoms at t_i and at residence

374 location x_i being infected by a case showing symptoms at time t_j and at residence location x_j is

$$f(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = \alpha(t_i - t_j - \gamma)e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)} e^{-\beta(x_i - x_j)^2} \quad (7)$$

375 The survival term is again determined by integrating the likelihood over all potential infection times

376 and all distances

$$S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = \left(\int_0^\infty \int_0^{t_i - t_j} \alpha(t_i - t_j - \gamma)e^{-\frac{1}{2}\alpha(t_i - t_j - \gamma)} e^{-\beta(x_i - x_j)^2} dt dx \right) \quad (8)$$

377 Integrating over time returns

$$S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = e^{-\frac{1}{2}\alpha(t_i-t_j-\gamma)} \int_0^{\infty} e^{-\beta(x_i-x_j)^2} dx \quad (9)$$

378 Integrating over all distances gives

$$S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = e^{-\frac{1}{2}\alpha(t_i-t_j-\gamma)} \frac{\sqrt{\pi}}{2\sqrt{\beta}} \quad (10)$$

379 Following equation 10, the hazard is equivalent to

$$H(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = \frac{\alpha(t_i - t_j - \gamma) e^{-\frac{1}{2}\alpha(t_i-t_j-\gamma)} e^{-\beta(x_i-x_j)^2}}{e^{-\frac{1}{2}\alpha(t_i-t_j-\gamma)} \frac{\sqrt{\pi}}{2\sqrt{\beta}}} \quad (11)$$

380 Which simplifies to

$$H(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) = \frac{2\sqrt{\beta}\alpha(t_i - t_j - \gamma) e^{-\beta(x_i-x_j)^2}}{\sqrt{\pi}} \quad (12)$$

381 Modelling missing cases using ϵ edges

382 The vast majority of disease surveillance and outbreak response datasets will not be able to capture
 383 all cases due to asymptomatic infection, underreporting and movement of people in/out of the
 384 surveillance area. Therefore, it is important to consider the impact of missing information on results
 385 and identify potential missing sources of infection. We use Epsilon edges, ϵ_i , to identify potential
 386 sources of infection. Here, each hazard is estimated as a further competing edge of transmission
 387 from an unobserved source, $H_0(\epsilon_i)$. Depending on assumptions for the likelihood and extent of
 388 unobserved infection sources, the epsilon edge value can be set to a high or low value. When high,
 389 we assume high amounts of unobserved infection and unless two cases have a very high likelihood
 390 of being linked, we assume the case was from an unobserved source. When low, we assume little
 391 missing data and so cases are only linked to an outside source if they are very unlikely to be linked to
 392 an observed candidate infector.

393 Adding epsilon, ϵ , as a competing hazard and survival returns:

394
$$f(\mathbf{t}, \mathbf{x}; \boldsymbol{\alpha}, \boldsymbol{\epsilon}, \boldsymbol{\beta}) =$$

395
$$\prod_{t_i \in \mathbf{t}} S_0(\epsilon_i) \prod_{I_k: t_k < t_i} S(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) \left(H_0(\epsilon_i) + \sum_{I_j: t_j < t_i} H(x_i, t_i | x_j, t_j; \alpha_{i,j}, \beta) \right) \quad (13)$$

396 The objective function is then:

397
$$\text{minimize}_{\boldsymbol{\alpha}, \boldsymbol{\epsilon}} -\log f(\mathbf{t}, \mathbf{x}; \boldsymbol{\alpha}, \boldsymbol{\epsilon}, \boldsymbol{\beta}) \quad \text{subject to } \boldsymbol{\alpha}, \boldsymbol{\epsilon}, \boldsymbol{\beta} > 0 \quad \forall i, j \quad (14)$$

398 Because this was carried out within a Bayesian framework the log posterior was maximised to obtain
 399 the maximum-a-posteriori estimates.

400 The algorithm was written in TensorFlow, implemented in R via the *rTensorflow* package. A prior
 401 probability was defined for the parameter shaping the serial interval of malaria, informed by
 402 previous characterisations of the serial interval of malaria²⁰. Because data about how likely the
 403 cases were to have moved long distances or the likelihood of a case has been infected by an
 404 unobserved source of infection were not available for the contexts explored here, several different
 405 parameterisations of the model were used to represent different scenarios (Table 1) and a detailed
 406 sensitivity analysis was carried out (Table 3). The versions of the model which are described in **Table**
 407 **1** and **Figures 13-15** represent different patterns of human/parasite movement, ranging from a
 408 context where there may be small amounts of movement (almost all under 10km) to moderate
 409 amounts of movement/travel(almost all under 50km) to a less restrictive parameterisation, where
 410 near cases were more likely but far away cases were not completely excluded. We applied different
 411 versions of the algorithm, as well as temporal-only algorithm to these datasets to explore the impact
 412 of different assumptions about the impact of space on estimated R_c values and their variation over
 413 time and space. We also evaluated the performance of each approach by comparing differences in
 414 the second order AIC (ΔAIC_c), and the corresponding Akaike Weights.

415 Twelve scenarios (**Table 1**) were considered when defining parameters for each dataset. These
 416 scenarios consider three different levels of likelihood of transmission in relationship to Euclidian
 417 distance (due to the limited range of mosquito travel, this is considered in the context of human

418 mobility), which was defined for both exponential and Gaussian kernels. These are illustrated in
419 **Figures 13-15**. Then the values for epsilon were set at 0.001 and 0.1, representing different levels of
420 missing cases likely. This can be interpreted as the chance of a case having an unobserved source of
421 infection. For example, 0.1 would represent $P(\text{unobserved source of infection}) = 0.1$.

422 The timeseries of R_c and its spatial patterns were illustrated for each dataset and parameter
423 combination and compared to the results of the time-only version of the algorithm. The results were
424 also mapped to compare how spatial patterns in R_c were affected by assumptions about space and
425 unobserved infections.

426 In order to compare models quantifiably, the second order Akaike Information Criterion (AICc) was
427 calculated using the equation $AIC_c = -2 \log f(x) + 2K(\frac{n}{n-k-1})$, where $f(x)$ is the model
428 likelihood, K is the number of parameters estimated and n is the sample size of the data used to fit
429 the parameters. The AIC²¹ is used in model comparison, by creating a comparison of negative log
430 likelihood that penalises increases in model parameters, to prevent overfitting. AICc is
431 recommended for use with smaller datasets with larger numbers of parameters, and as the sample
432 size n increases AICc converges to AIC¹⁰. The differences in AICc for each model, known as $\Delta AICc$,
433 were calculated to compare models. Typically, a $\Delta AICc$ of greater than 10 is considered strong
434 evidence that that model performs worse than the model it is being compared to.

435 In addition, Akaike Weights were calculated, which are a measure of the relative likelihood of a
436 model compared to the others considered. Akaike weights are determined by taking the normalised
437 relative likelihood of a model which is $\exp(-0.5 * \Delta AICc \text{ score})$, and then dividing by the sum of
438 these values across all models to obtain a normalised result.

439 **Sensitivity analysis and comparison of prior choice on estimated results**

440 In the scenario analysis above the distance shaping parameter is fixed. However due to the
441 uncertainties in the relationship between distance and likelihood of transmission, in many contexts

442 it may be useful to estimate β . To explore the relationship between the estimated epsilon edges, ϵ ,
443 and estimated shaping parameter, β , for the distance function. a detailed sensitivity analysis was
444 carried out to explore the impact of a) prior choice for ϵ d) prior choice for β on both the maximum-
445 a-posteriori estimates for β and the estimated mean R_c .

446 To consider the effect of varying parameter values and explore their interactions, a range of distance
447 and epsilon edge priors were considered. A truncated normal prior was used for both parameters,
448 and the mean and standard deviation were varied. For ϵ the mean was varied between 1e-10 and
449 0.5, and the standard deviation was varied between 0.0001 and 0.1. For β , the mean for a Gaussian
450 Kernel was varied between 0.00001 and 0.01 and for an exponential kernel the means considered
451 ranged between 0.0001 and 0.1. For both the standard deviations varied between 0.0001 and 0.1
452 (Table). Every possible combination of the parameters were run for each dataset and both Gaussian
453 and exponential spatial kernels, giving a total of 2400 parameter combinations tested per kernel, per
454 dataset.

455

456

457

[Insert Figures 13-15 approximately here]

458

459

460 Bibliography

- 461 1. Lourenço, C. *et al.* Strengthening surveillance systems for malaria elimination: a global
462 landscaping of system performance, 2015–2017. *Malar. J.* **18**, 1–11 (2019).
- 463 2. Sturrock, H. J. W. *et al.* Mapping Malaria Risk in Low Transmission Settings: Challenges and
464 Opportunities. *Trends Parasitol.* **32**, 635–645 (2016).
- 465 3. Wesolowski, A. *et al.* Mapping malaria by combining parasite genomic and epidemiologic data.
466 *BMC Med.* **16**, 190 (2018).
- 467 4. Prothero, R. M. Disease and Mobility: A Neglected Factor in Epidemiology. *Int. J. Epidemiol.* **6**,
468 259–267 (1977).
- 469 5. Pindolia, D. K. *et al.* Human movement data for malaria control and elimination strategic
470 planning. *Malar. J.* **11**, 205 (2012).
- 471 6. Buckee, C. O., Wesolowski, A., Eagle, N. N., Hansen, E. & Snow, R. W. Mobile phones and
472 malaria: Modeling human and parasite travel. *Travel Med. Infect. Dis.* **11**, 15–22 (2013).
- 473 7. Wesolowski, A. *et al.* Quantifying the Impact of Human Mobility on Malaria. *Science* **338**, 267–
474 270 (2012).
- 475 8. Cotter, C. *et al.* The changing epidemiology of malaria elimination: new strategies for new
476 challenges. *The Lancet* **382**, 900–911 (2013).
- 477 9. Routledge, I. *et al.* Tracking progress towards malaria elimination in China: Individual-level
478 estimates of transmission and its spatiotemporal variation using a diffusion network approach.
479 *PLOS Comput. Biol.* **16**, e1007707 (2020).
- 480 10. Hurvich, C. M. & Tsai, C. L. Regression and time series model selection in small samples.
481 *Biometrika* **76**, 297–307 (1989).
- 482 11. Routledge, I. *et al.* Estimating spatiotemporally varying malaria reproduction numbers in a near
483 elimination setting. *Nat. Commun.* **9**, 2476 (2018).
- 484 12. Bateman, A. J. Is gene dispersion normal? *Heredity* **4**, 353–363 (1950).

- 485 13. Wang, L., Ermon, S. & Hopcroft, J. E. Feature-Enhanced Probabilistic Models for Diffusion
486 Network Inference. in *Lecture Notes in Computer Science (including subseries Lecture Notes in*
487 *Artificial Intelligence and Lecture Notes in Bioinformatics)* vol. 7524 LNAI 499–514 (2012).
- 488 14. Simini, F., González, M. C., Maritan, A. & Barabási, A.-L. A universal model for mobility and
489 migration patterns. *Nature* **484**, 96–100 (2012).
- 490 15. Weiss, D. J. *et al.* A global map of travel time to cities to assess inequalities in accessibility in
491 2015. *Nature* **553**, 333–336 (2018).
- 492 16. Bousema, T. & Drakeley, C. Epidemiology and Infectivity of *Plasmodium falciparum* and
493 *Plasmodium vivax* Gametocytes in Relation to Malaria Control and Elimination. *Clin. Microbiol.*
494 *Rev.* **24**, 377–410 (2011).
- 495 17. Sturrock, H. J. W. *et al.* Mapping Malaria Risk in Low Transmission Settings: Challenges and
496 Opportunities. *Trends in Parasitology* vol. 32 635–645 (2016).
- 497 18. Reiner, R. C. *et al.* Mapping residual transmission for malaria elimination. *eLife* **4**, e09520 (2015).
- 498 19. Lai, S. *et al.* Changing epidemiology and challenges of malaria in China towards elimination.
499 *Malar. J.* **18**, 107 (2019).
- 500 20. Huber, J. H., Johnston, G. L., Greenhouse, B., Smith, D. L. & Perkins, T. A. Quantitative, model-
501 based estimates of variability in the generation and serial intervals of *Plasmodium falciparum*
502 malaria. *Malar. J.* **15**, 490 (2016).
- 503 21. Akaike, H. A New Look at the Statistical Model Identification. *Autom. Control IEEE Trans. On* **19**,
504 716–723 (1974).

505
506
507
508
509
510
511

512 Author contributions

513 IR devised the project, carried out analysis, curated data, wrote paper. SB supervised the project and
514 edited paper draft.

515 Additional Information

516 The authors declare no competing interests.

517 Code for algorithm to generate reproduction numbers, estimate unobserved cases and calculate
518 metrics of model performance are available at https://github.com/IzzyRou/spatial_rcs

519 Supplementary Information and Figure Legends

520 Supplementary Information

521 **Supplementary Table 1:** Full results of ΔAICc and Akaike Weights for each scenario, dataset and
522 spatial kernel considered

523

524 Figure Legends

525 **Figure 1:** R_c estimates from El Salvador line list based on using the time-only scenario and Scenarios
526 1-12 with an exponential kernel

527 **Figure 2 :** R_c estimates from Eswatini line list based on using the time-only scenario and Scenarios 1-
528 12 with an exponential kernel

529 **Figure 3:** R_c estimates from China *P. falciparum* line list based on using the time-only scenario and
530 Scenarios 1-12 with an exponential kernel

531 **Figure 4:** R_c estimates from China *P. vivax* line list based on using the time-only scenario and
532 Scenarios 1-12 with an exponential kernel

533 **Figure 5: Map of R_c estimates for El Salvador** Map of A) Time-only B) Best scenario by AIC (Scenario
534 9) and C) Scenario 4, representing an assumption of little long-distance transmission and few
535 unobserved cases. Note the increasing focality in C), with higher R_c values estimated on the Pacific
536 Coastal area of the Ahuacapan and Sonsonote municipalities, where the NMCP have long identified
537 as the remaining foci of risk.

538 **Figure 6: Map of R_c estimates for Swaziland** Map of A) Time-only B) Best scenario by AIC (Scenario
539 9) and C) Scenario 4, representing an assumption of little long-distance transmission and few
540 unobserved cases. Note the increasing focality in C), with higher R_c values estimated around the
541 northern corner of the country which borders Mozambique.

542 **Figure 7: Map of R_c estimates for *P. falciparum* in China** Map of A) Time-only B) Best scenario by
543 AIC (Scenario 11) and C) Scenario 4, representing an assumption of little long-distance transmission
544 and few unobserved cases.

545 **Figure 8: Map of R_c estimates for *P. vivax* in China** Map of A) Time-only B) Best scenario by AIC
546 (Scenario 11) and C) Scenario 4, representing an assumption of little long-distance transmission and
547 few unobserved cases.

548 **Figure 9: El Salvador sensitivity analysis.** Sensitivity analysis showing the impact of varying the prior
549 mean for the distance kernel shaping parameter, β . The different colours and shapes represent

550 different means and standard deviations respectively of the normally-distributed prior of epsilon,
551 ϵ , which represents shapes represent different hazards of infection by an external, unobserved
552 source. For A-D, the x-axis represents the prior mean used for β . A) the y-axis shows the maximum a
553 posteriori parameter estimate for the parameter β . B) shows the same results, stratified by the prior
554 mean of ϵ for clarity. C) Shows the impact of priors for β and ϵ on the mean R_c estimate, and again
555 D) shows the same result, stratified by the prior mean of ϵ .

556 **Figure 10: Eswatini sensitivity analysis.** Sensitivity analysis showing the impact of varying the prior
557 means for Eswatini. Sensitivity analysis showing the impact of varying the prior mean for the
558 distance kernel shaping parameter, β . The different colours and shapes represent different means
559 and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which represents
560 shapes represent different hazards of infection by an external, unobserved source. For A-D, the x-
561 axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori parameter
562 estimate for the parameter β . B) shows the same results, stratified by the prior mean of ϵ for clarity.
563 C) Shows the impact of priors for β and ϵ on the mean R_c estimate, and again D) shows the same
564 result, stratified by the prior mean of ϵ .

565 **Figure 11: China (*P. falciparum*) Sensitivity Analysis.** Sensitivity analysis showing the impact of
566 varying the prior means for *P. falciparum* in China. Sensitivity analysis showing the impact of varying
567 the prior means for Eswatini. Sensitivity analysis showing the impact of varying the prior mean for
568 the distance kernel shaping parameter, β . The different colours and shapes represent different
569 means and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which
570 represents shapes represent different hazards of infection by an external, unobserved source. For A-
571 D, the x-axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori
572 parameter estimate for the parameter β . B) shows the same results, stratified by the prior mean of
573 ϵ for clarity. C) Shows the impact of priors for β and ϵ on the mean R_c estimate, and again D) shows
574 the same result, stratified by the prior mean of ϵ .

575 **Figure 12: China (*P. vivax*) Sensitivity Analysis.** Sensitivity analysis showing the impact of varying
576 the prior means for *P. vivax* in China. Sensitivity analysis showing the impact of varying the prior
577 means for Eswatini. Sensitivity analysis showing the impact of varying the prior mean for the
578 distance kernel shaping parameter, β . The different colours and shapes represent different
579 means and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which represents
580 shapes represent different hazards of infection by an external, unobserved source. For A-D, the x-
581 axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori parameter
582 estimate for the parameter β . B) shows the same results, stratified by the prior mean of ϵ for clarity.
583 C) Shows the impact of priors for β and ϵ on the mean R_c estimate, and again D) shows the same
584 result, stratified by the prior mean of ϵ .

585 **Figure 13: Illustration of likelihoods, hazards and survivals for less restrictive kernels (longer range
586 human movement likely).** Plots showing how the pairwise likelihoods, survivals and hazards vary
587 with time and distance under different model structures. The first row of plots shows the pairwise
588 likelihoods, the second row shows the pairwise survival and the third row shows the pairwise hazard
589 values for different combinations of distance (in kilometres) and time between symptom onset
590 (days). The first column shows the results for a time-only version of the algorithm. The second
591 column shows results for an exponential kernel and the third column shows results for a Gaussian
592 kernel. In this example less restrictive values for beta, the shaping parameter for the distance
593 kernels have been chosen, representing a context where there is more long-range movement of
594 parasites.

595 **Figure 14: Illustration of likelihoods, hazards and survivals for moderately restrictive kernels
596 (moderate human movement, most movement under 50km).** Plots showing how the pairwise
597 likelihoods, survivals and hazards vary with time and distance under different model structures. The

598 first row of plots shows the pairwise likelihoods, the second row shows the pairwise survival and the
599 third row shows the pairwise hazard values for different combinations of distance (in kilometres)
600 and time between symptom onset (days). The first column shows the results for a time-only version
601 of the algorithm. The second column shows results for an exponential kernel and the third column
602 shows results for a Gaussian kernel. In this example values for beta, the shaping parameter for the
603 distance kernels have been chosen to represent a context where there is more some movement of
604 parasites, but where little movement is expected beyond 50-75km. The likelihood for the Gaussian
605 Kernel is more concentrated, which could represent shorter range movement e.g. commutes,
606 whereas the Exponential has a longer tail so could represent a mixture of short and longer range
607 parasite movement.

608 **Figure 15: Illustration of likelihoods, hazards and survivals for highly restrictive kernels (Human**
609 **movement unlikely, most movement under 10km).** Plots showing how the pairwise likelihoods,
610 survivals and hazards vary with time and distance under different model structures. The first row of
611 plots shows the pairwise likelihoods, the second row shows the pairwise survival and the third row
612 shows the pairwise hazard values for different combinations of distance (in kilometres) and time
613 between symptom onset (days). The first column shows the results for a time-only version of the
614 algorithm. The second column shows results for an exponential kernel and the third column shows
615 results for a Gaussian kernel. In this example more restrictive values for beta, the shaping
616 parameter for the distance kernels have been chosen, representing a context where there is very
617 little movement of parasites, with very little movement beyond 10-20km expected.

618

Figures

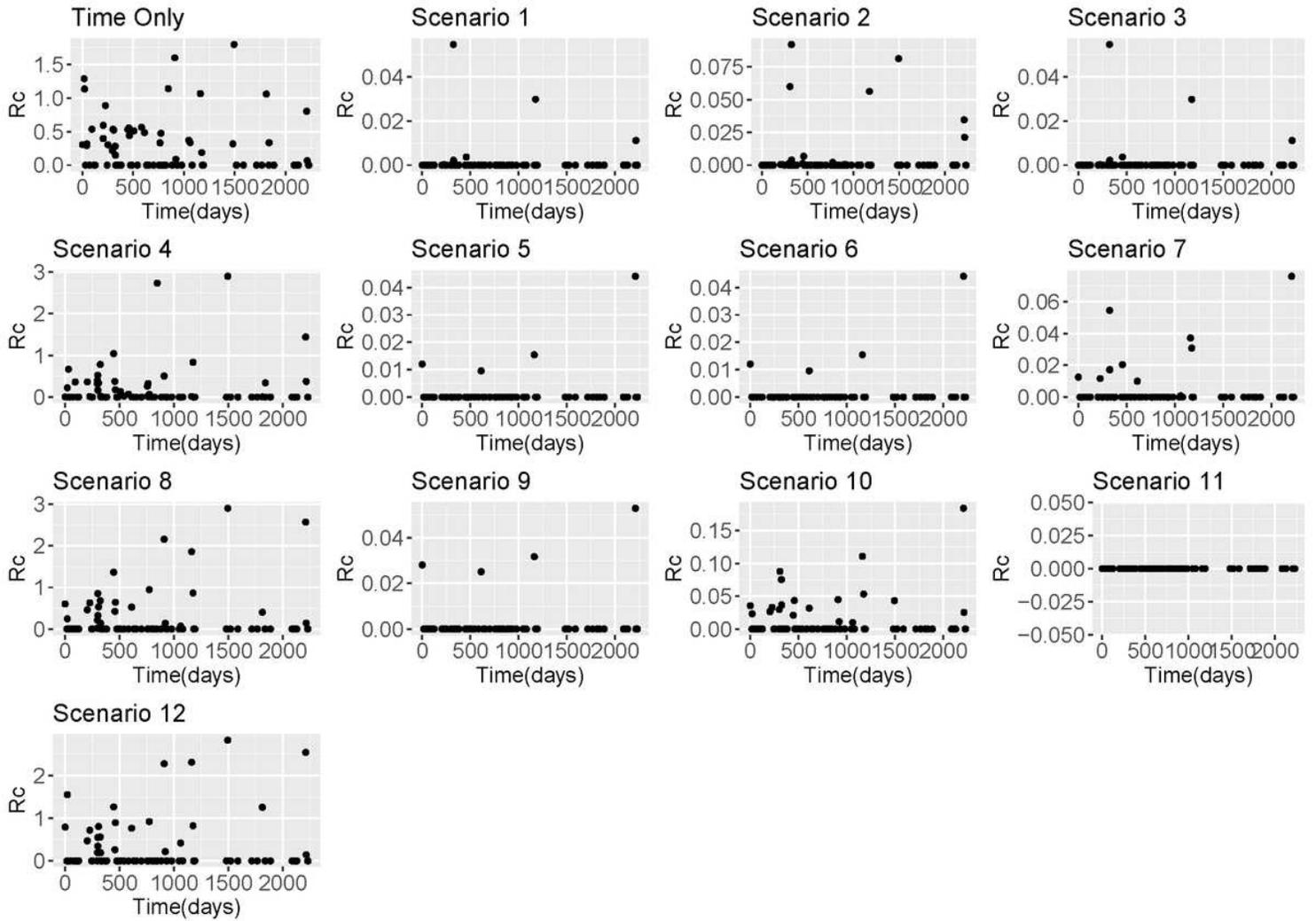


Figure 1

R_c estimates from El Salvador line list based on using the time-only scenario and Scenarios 1-12 with an exponential kernel

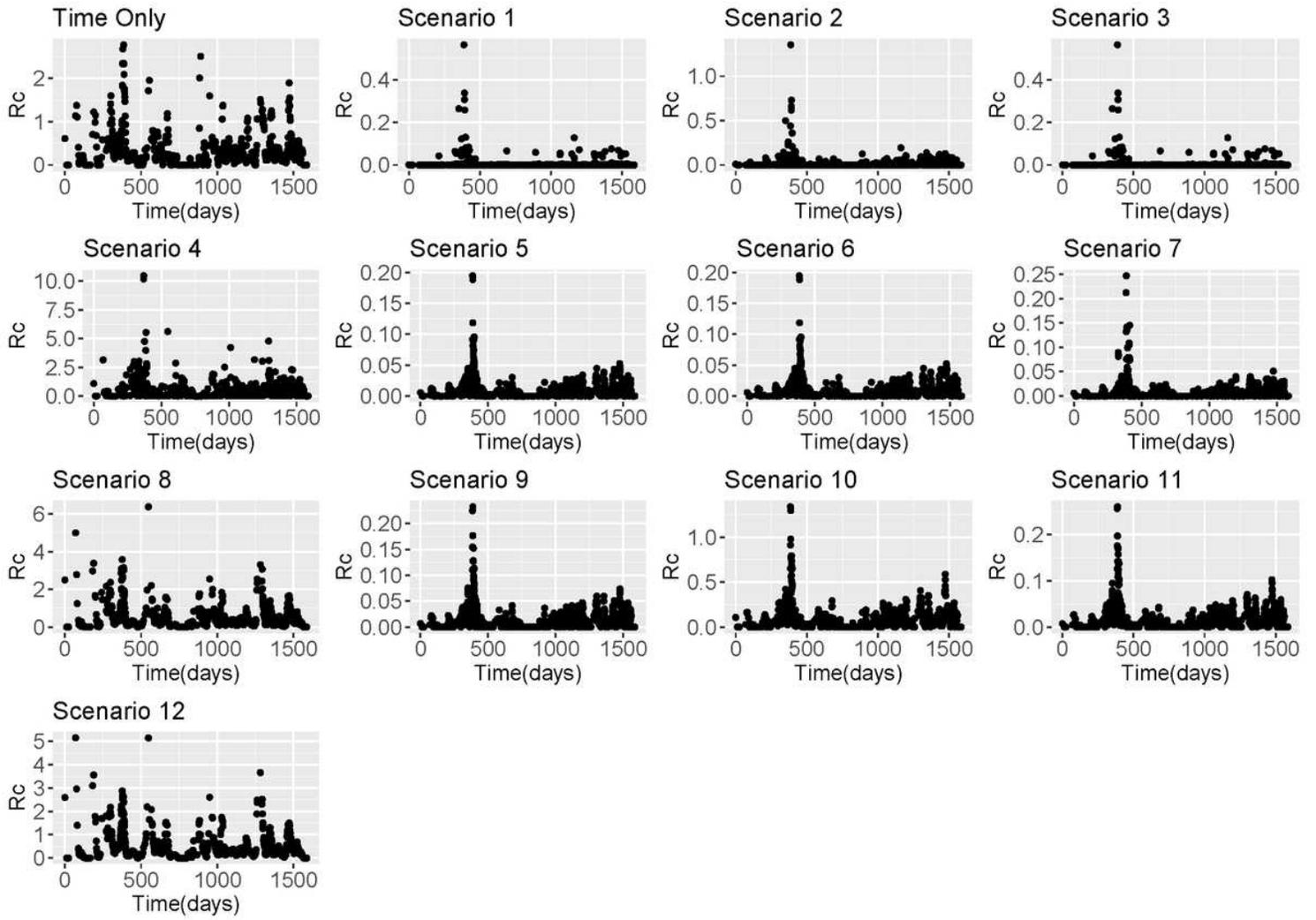


Figure 2

R_c estimates from Eswatini line list based on using the time-only scenario and Scenarios 1-12 with an exponential kernel

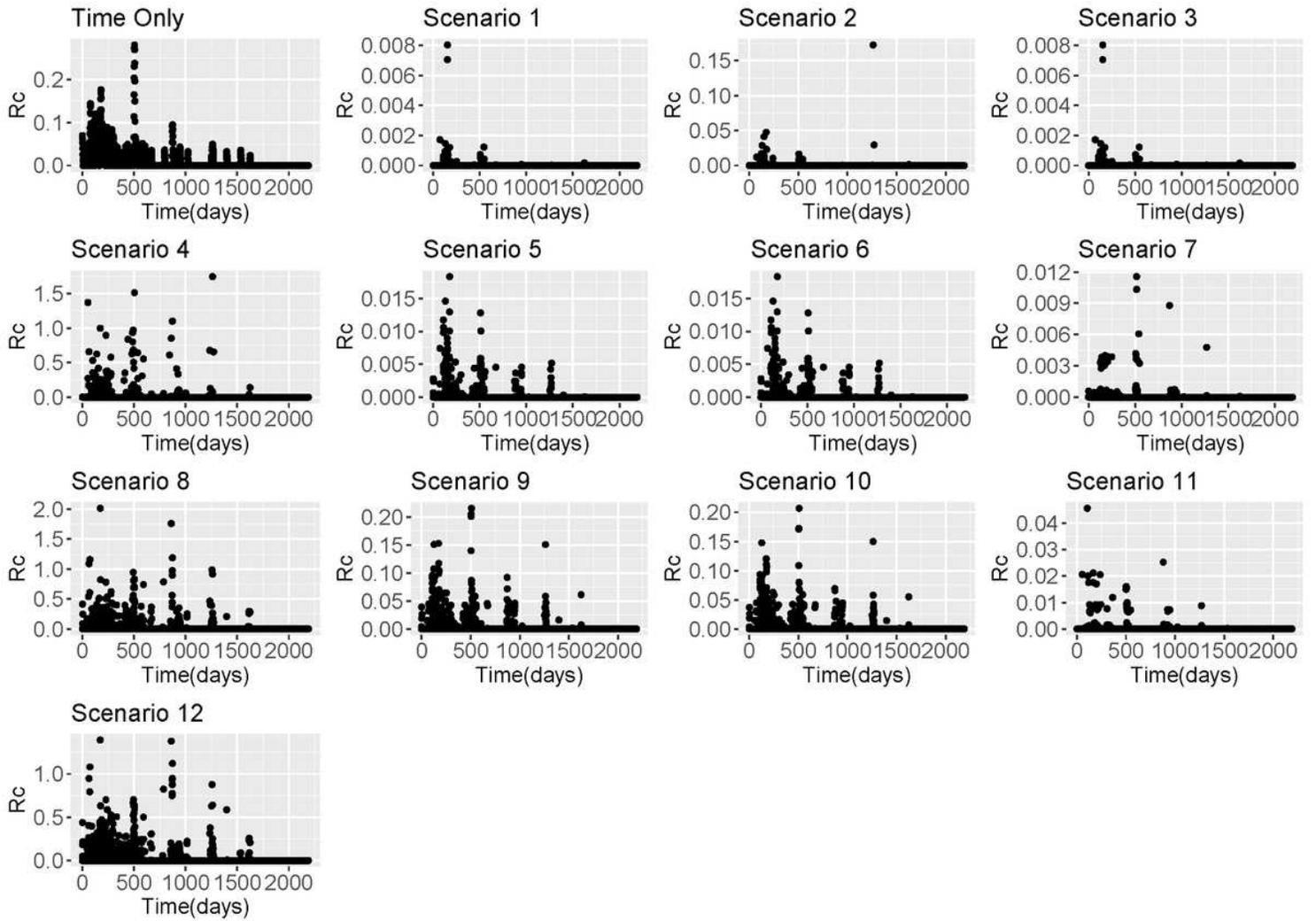


Figure 3

Rcestimates from China *P. falciparum* line list based on using the time-only scenario and Scenarios 1-12 with an exponential kernel

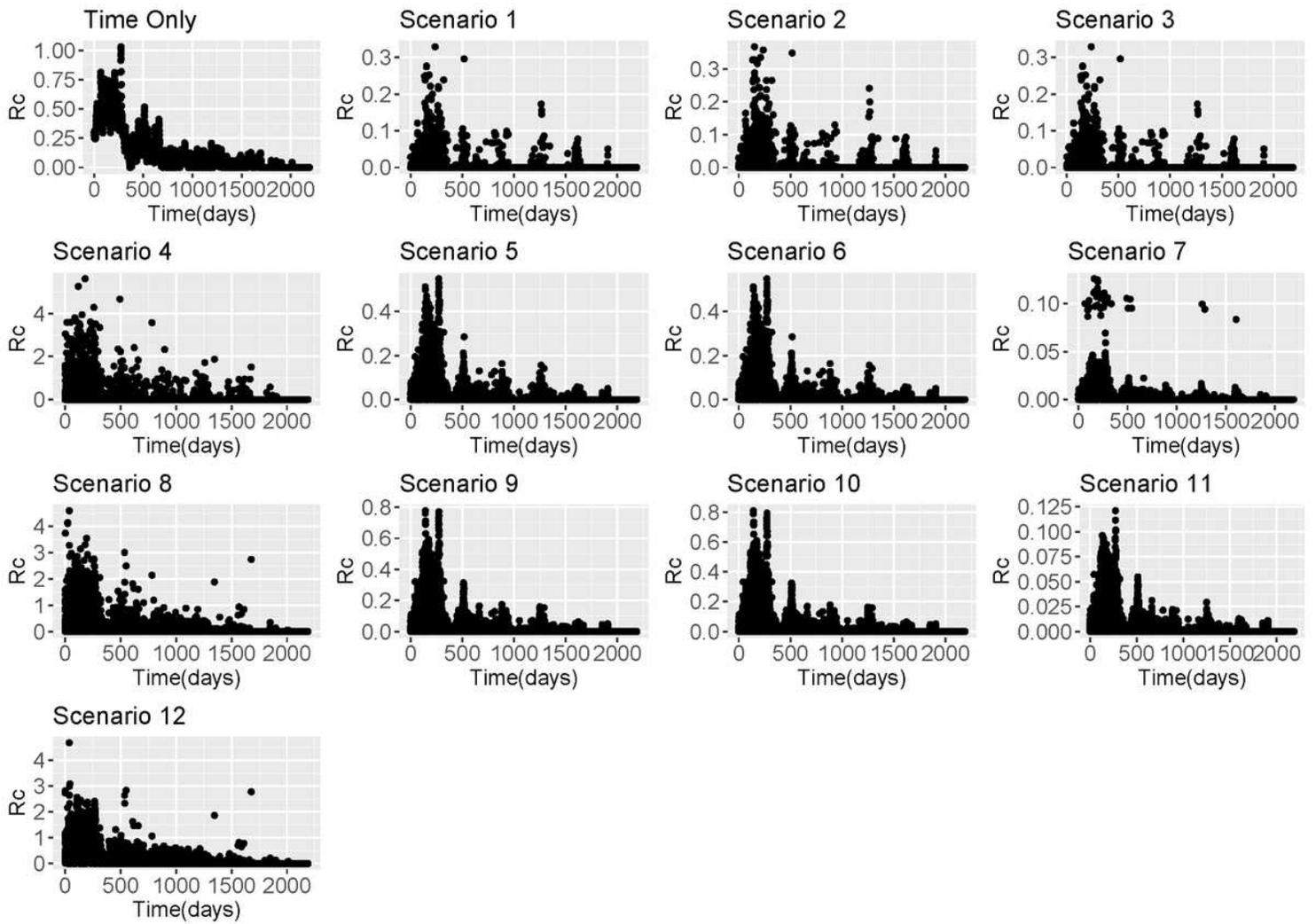


Figure 4

Rcestimates from China *P. vivax* line list based on using the time-only scenario and Scenarios 1-12 with an exponential kernel

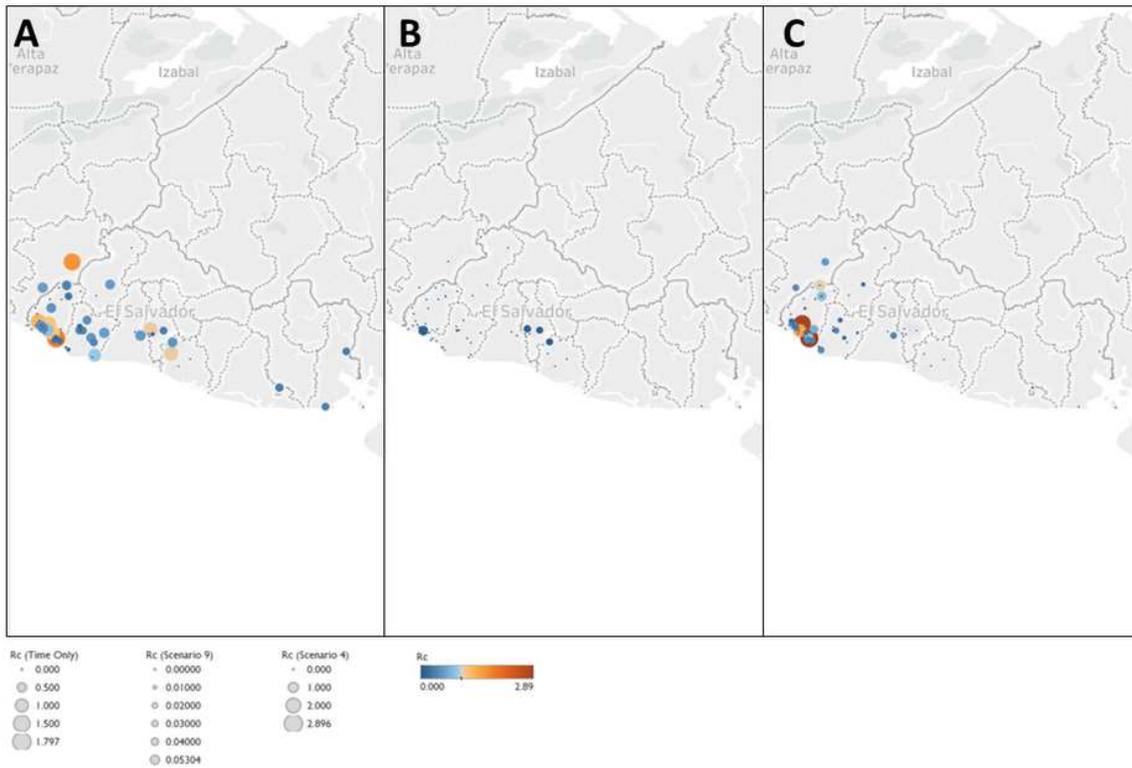


Figure 5

Map of R_c estimates for El Salvador Map of A) Time-only B) Best scenario by AIC (Scenario 9) and C) Scenario 4, representing an assumption of little long-distance transmission and few unobserved cases. Note the increasing focality in C), with higher R_c values estimated on the Pacific Coastal area of the Ahuacapan and Sonsonote municipalities, where the NMCP have long identified as the remaining foci of risk.

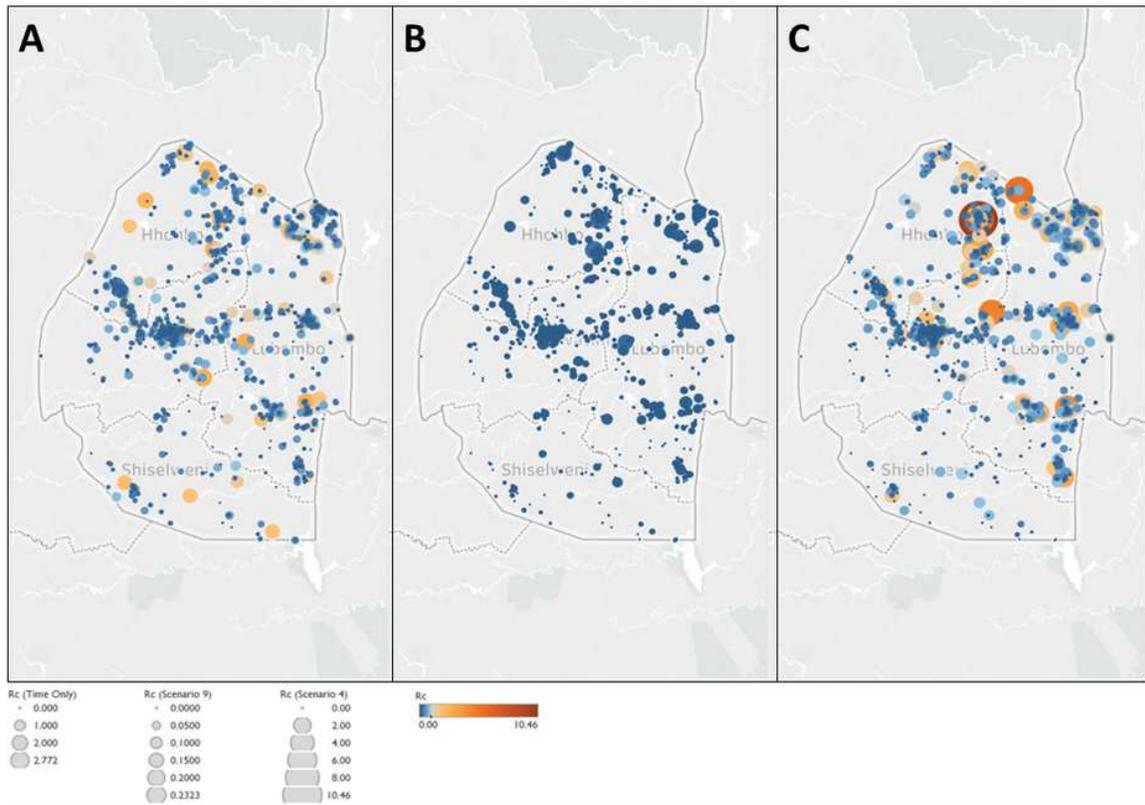


Figure 6

Map of R_c estimates for Swaziland Map of A) Time-only B) Best scenario by AIC (Scenario 9) and C) Scenario 4, representing an assumption of little long-distance transmission and few unobserved cases. Note the increasing focality in C), with higher R_c values estimated around the northern corner of the country which borders Mozambique.

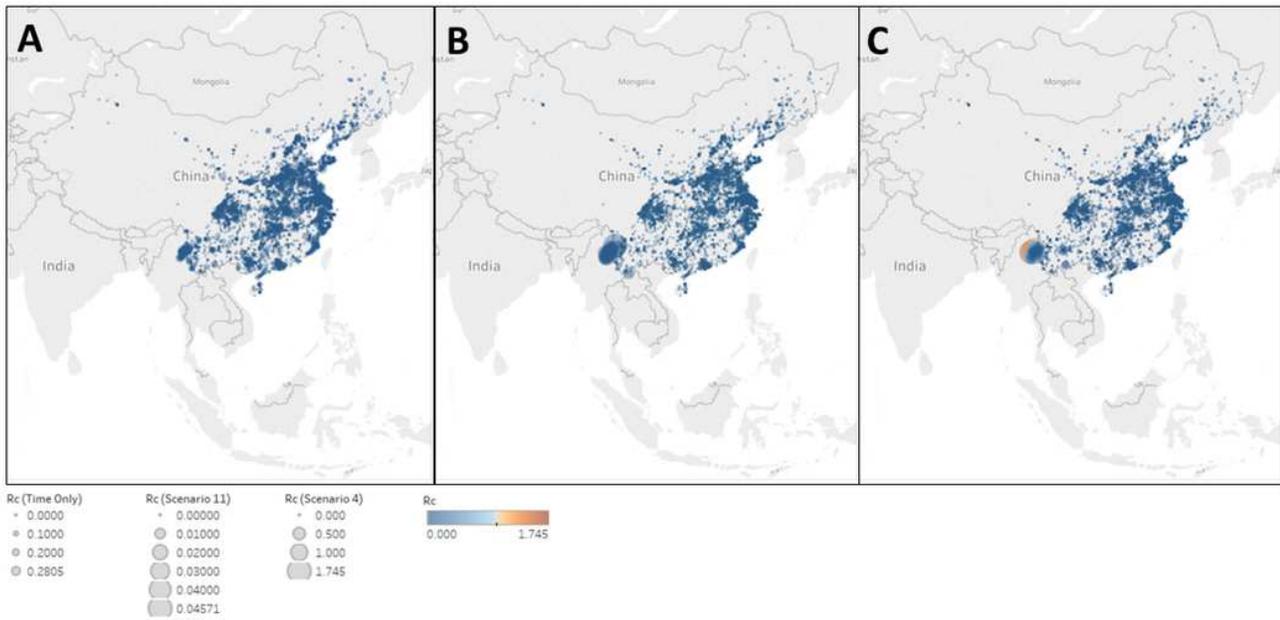


Figure 7

Map of R_c estimates for *P. falciparum* in China Map of A) Time-only B) Best scenario by AIC (Scenario 11) and C) Scenario 4, representing an assumption of little long-distance transmission and few unobserved cases.

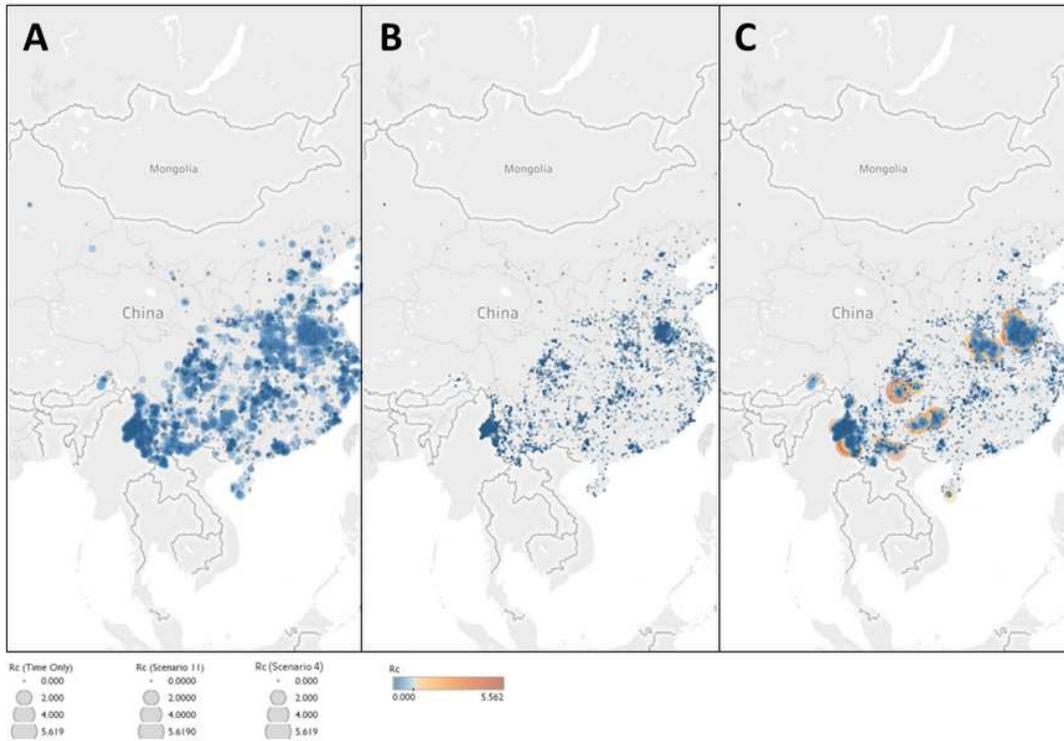


Figure 8

Map of R_c estimates for *P. vivax* in China Map of A) Time-only B) Best scenario by AIC (Scenario 11) and C) Scenario 4, representing an assumption of little long-distance transmission and few unobserved cases.

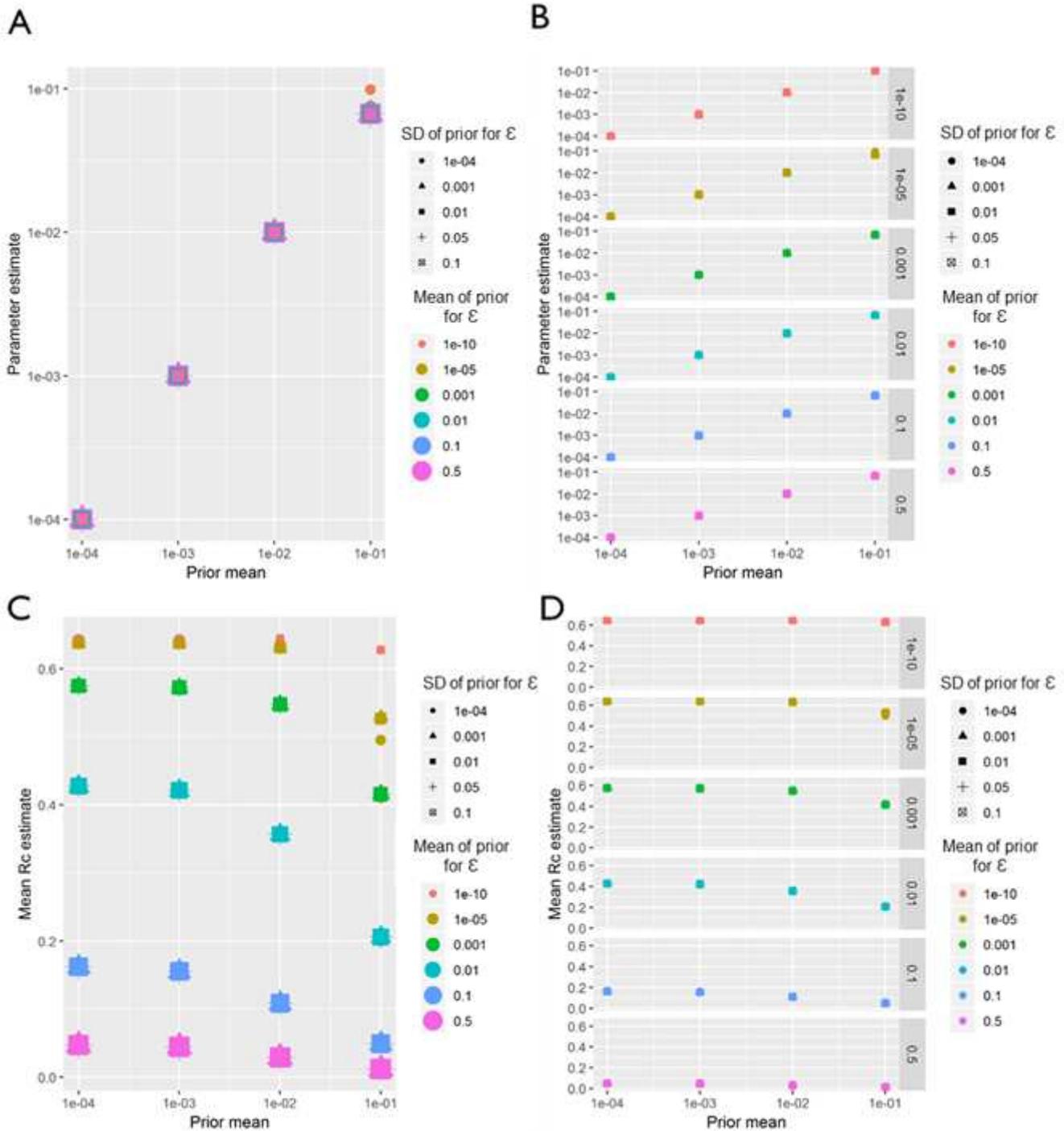


Figure 9

El Salvador sensitivity analysis. Sensitivity analysis showing the impact of varying the prior mean for the distance kernel shaping parameter, β . The different colours and shapes represent different means and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which represents shapes represent different hazards of infection by an external, unobserved source. For A-D, the x-axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori parameter estimate for the

parameter β . B) shows the same results, stratified by the prior mean of β for clarity. C) Shows the impact of priors for β and β on the mean R_c estimate, and again D) shows the same result, stratified by the prior mean of β .

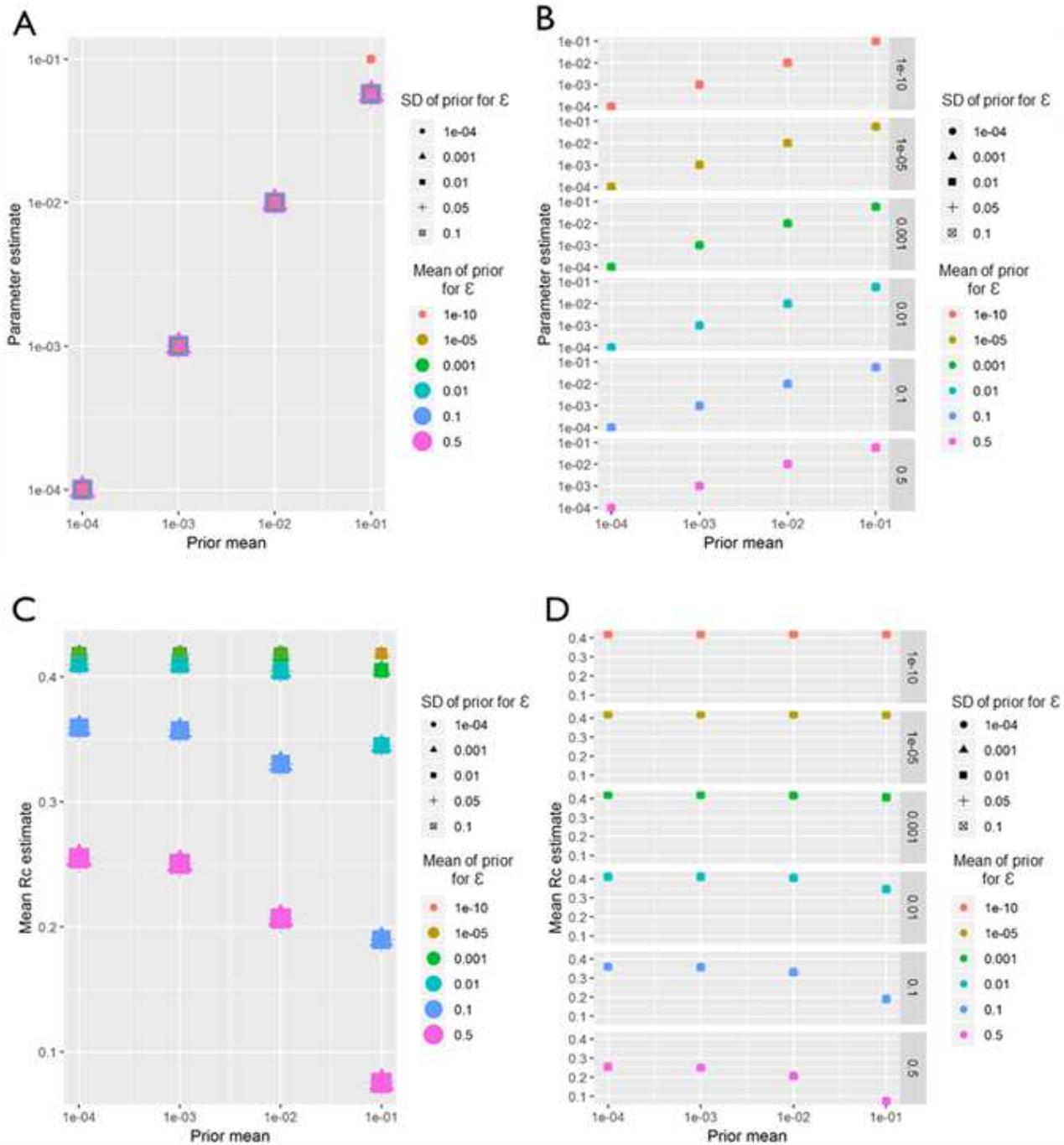


Figure 10

Eswatini sensitivity analysis. Sensitivity analysis showing the impact of varying the prior means for Eswatini. Sensitivity analysis showing the impact of varying the prior mean for the distance kernel

shaping parameter, β . The different colours and shapes represent different means and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which represents shapes represent different hazards of infection by an external, unobserved source. For A-D, the x-axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori parameter estimate for the parameter β . B) shows the same results, stratified by the prior mean of ϵ for clarity. C) Shows the impact of priors for β and ϵ on the mean R_c estimate, and again D) shows the same result, stratified by the prior mean of ϵ .

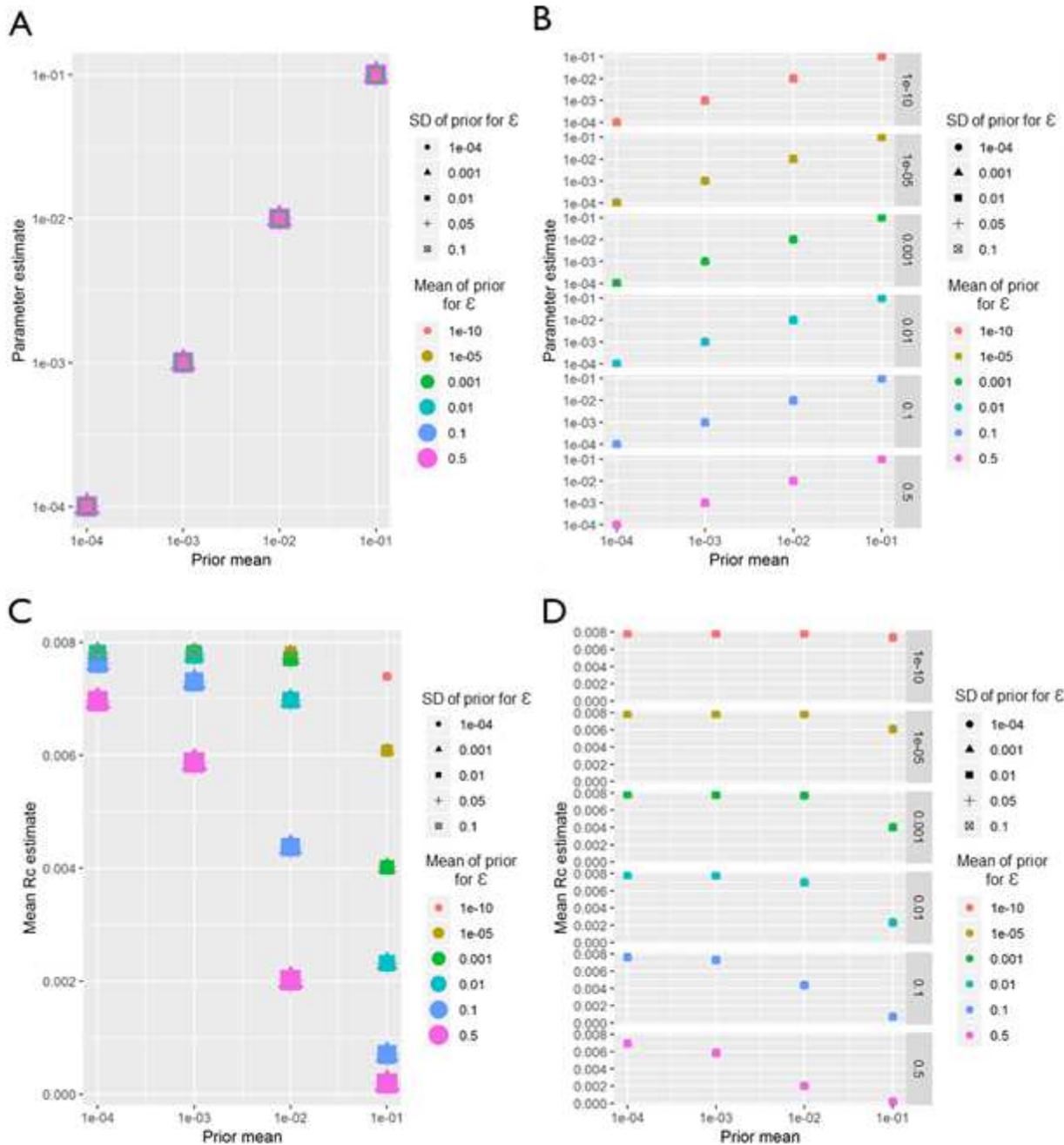


Figure 11

China (*P. falciparum*) Sensitivity Analysis. Sensitivity analysis showing the impact of varying the prior means for *P. falciparum* in China. Sensitivity analysis showing the impact of varying the prior means for Eswatini. Sensitivity analysis showing the impact of varying the prior mean for the distance kernel shaping parameter, β . The different colours and shapes represent different means and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which represents shapes represent different hazards of infection by an external, unobserved source. For A-D, the x-axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori parameter estimate for the parameter β . B) shows the same results, stratified by the prior mean of ϵ for clarity. C) Shows the impact of priors for β and ϵ on the mean R_c estimate, and again D) shows the same result, stratified by the prior mean of ϵ .

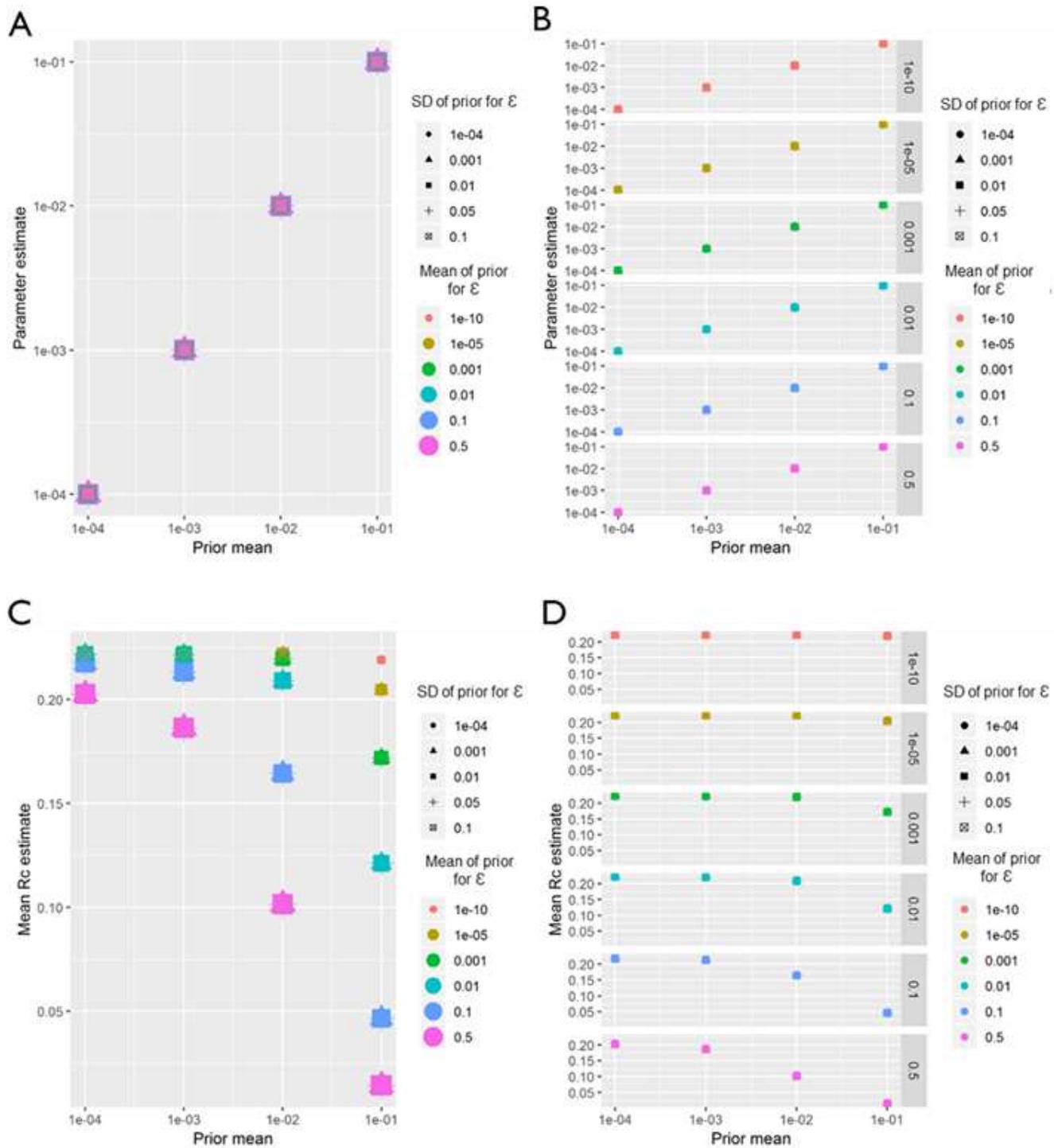


Figure 12

China (*P. vivax*) Sensitivity Analysis. Sensitivity analysis showing the impact of varying the prior means for *P. vivax* in China. Sensitivity analysis showing the impact of varying the prior means for Eswatini. Sensitivity analysis showing the impact of varying the prior mean for the distance kernel shaping parameter, β . The different colours and shapes represent different means and standard deviations respectively of the normally-distributed prior of epsilon, ϵ , which represents shapes represent different

hazards of infection by an external, unobserved source. For A-D, the x-axis represents the prior mean used for β . A) the y-axis shows the maximum a posteriori parameter estimate for the parameter β . B) shows the same results, stratified by the prior mean of β for clarity. C) Shows the impact of priors for β and β_0 on the mean R_c estimate, and again D) shows the same result, stratified by the prior mean of β .

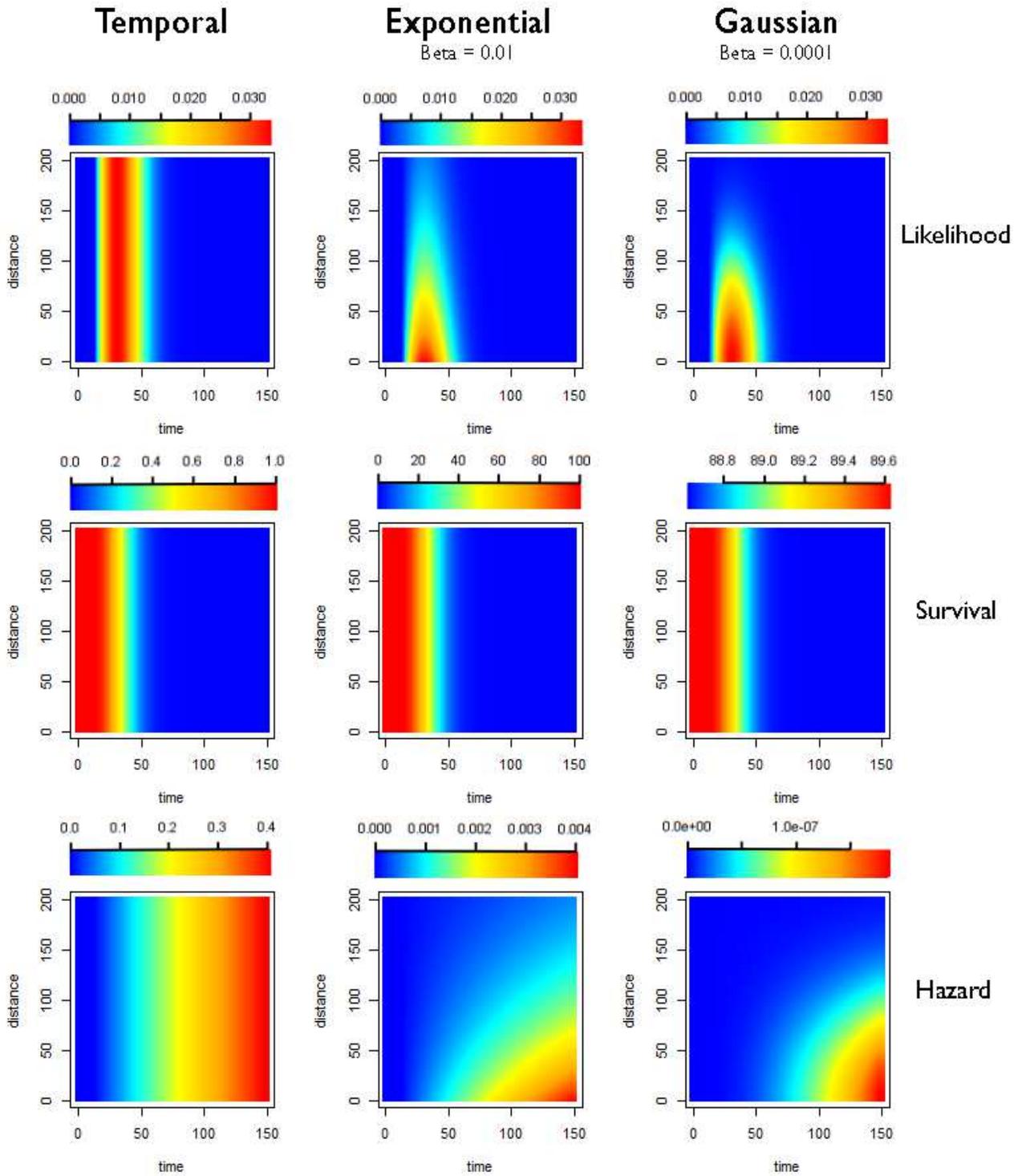


Figure 13

Illustration of likelihoods, hazards and survivals for less restrictive kernels (longer range human movement likely). Plots showing how the pairwise likelihoods, survivals and hazards vary with time and distance under different model structures. The first row of plots shows the pairwise likelihoods, the second row shows the pairwise survival and the third row shows the pairwise hazard values for different combinations of distance (in kilometres) and time between symptom onset (days). The first column shows the results for a time-only version of the algorithm. The second column shows results for an exponential kernel and the third column shows results for a Gaussian kernel. In this example less restrictive values for beta, the shaping parameter for the distance kernels have been chosen, representing a context where there is more long-range movement of parasites.

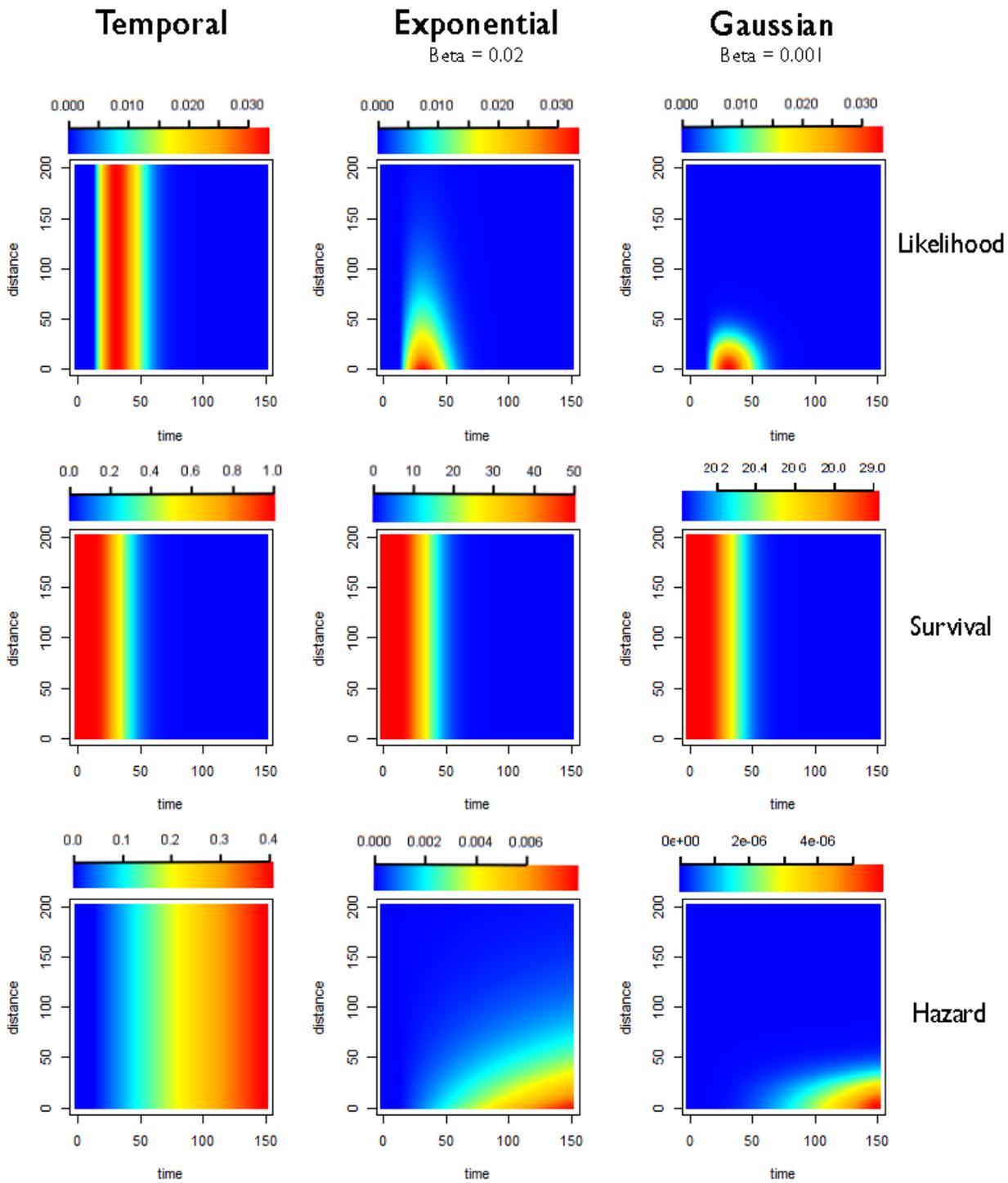


Figure 14

Illustration of likelihoods, hazards and survivals for moderately restrictive kernels (moderate human movement, most movement under 50km). Plots showing how the pairwise likelihoods, survivals and hazards vary with time and distance under different model structures. The first row of plots shows the pairwise likelihoods, the second row shows the pairwise survival and the third row shows the pairwise hazard values for different combinations of distance (in kilometres) and time between symptom onset

(days). The first column shows the results for a time-only version of the algorithm. The second column shows results for an exponential kernel and the third column shows results for a Gaussian kernel. In this example values for beta, the shaping parameter for the distance kernels have been chosen to represent a context where there is more some movement of parasites, but where little movement is expected beyond 50-75km. The likelihood for the Gaussian Kernel is more concentrated, which could represent shorter range movement e.g. commutes, whereas the Exponential has a longer tail so could represent a mixture of short and longer range parasite movement.

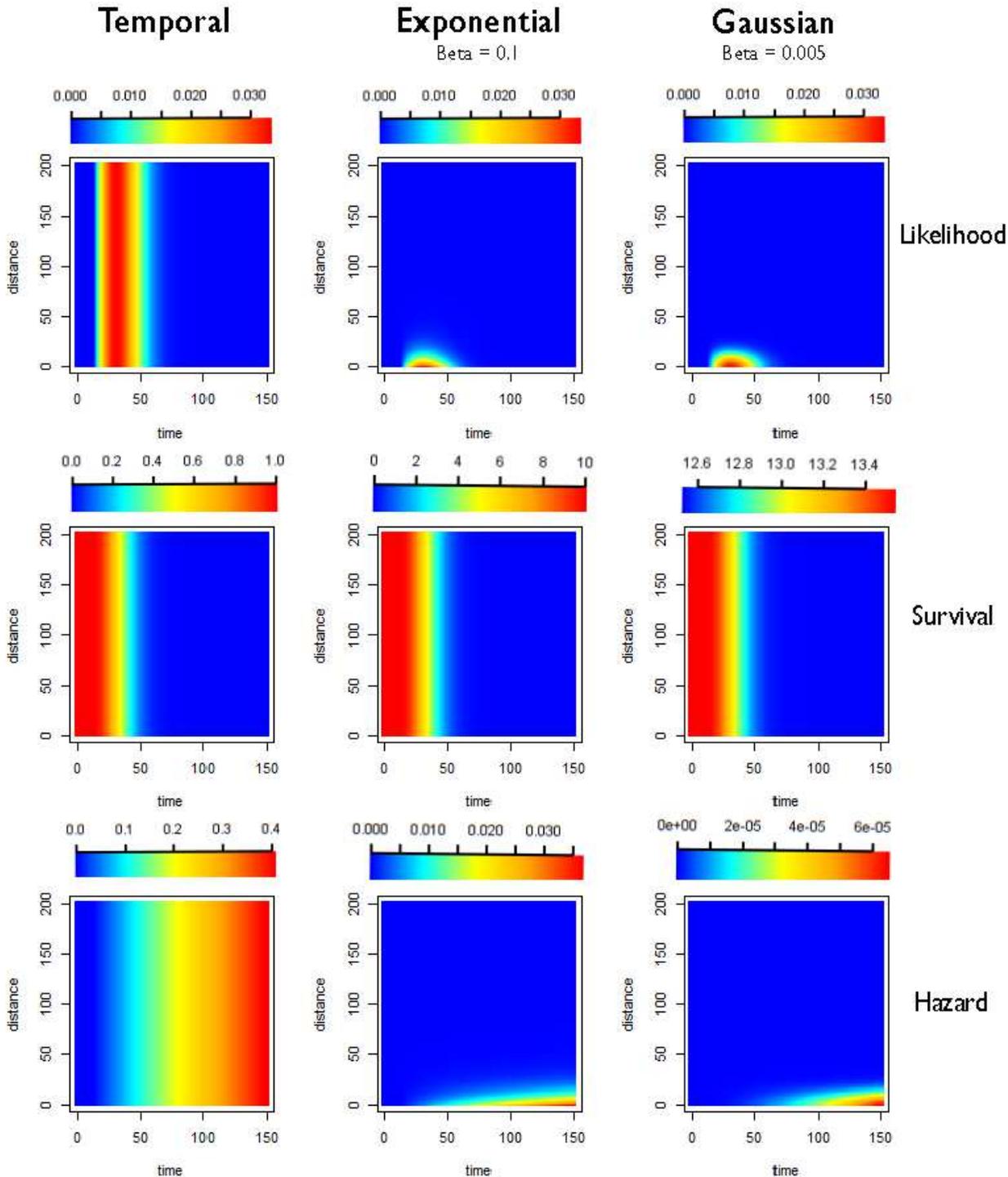


Figure 15

Illustration of likelihoods, hazards and survivals for highly restrictive kernels (Human movement unlikely, most movement under 10km). Plots showing how the pairwise likelihoods, survivals and hazards vary with time and distance under different model structures. The first row of plots shows the pairwise likelihoods, the second row shows the pairwise survival and the third row shows the pairwise hazard values for different combinations of distance (in kilometres) and time between symptom onset (days). The first column shows the results for a time-only version of the algorithm. The second column shows results for an exponential kernel and the third column shows results for a Gaussian kernel. In this example more restrictive values for beta, the shaping parameter for the distance kernels have been chosen, representing a context where there is very little movement of parasites, with very little movement beyond 10-20km expected.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryTable1.xlsx](#)
- [Titlesheet.docx](#)