

Keywords

Approximate Bayesian computation, outbreak dynamics, stochastic birth-death process, tuberculosis.

Corresponding author: Jukka Corander (jukka.corander@medisin.uio.no)

Author roles: **Lintusaari J:** Formal Analysis, Investigation, Methodology, Writing – Original Draft Preparation; **Blomstedt P:** Formal Analysis, Writing – Original Draft Preparation, Writing – Review & Editing; **Rose B:** Writing – Review & Editing; **Sivula T:** Formal Analysis, Writing – Original Draft Preparation, Writing – Review & Editing; **Gutmann MU:** Methodology, Writing – Original Draft Preparation, Writing – Review & Editing; **Kaski S:** Methodology, Supervision, Writing – Original Draft Preparation, Writing – Review & Editing; **Corander J:** Methodology, Supervision, Writing – Original Draft Preparation, Writing – Review & Editing

Competing interests: No competing interests were disclosed.

Grant information: This work was supported by the Academy of Finland (Finnish Centre of Excellence in Computational Inference Research, COIN; grants 294238 and 292334), the ERC (grant 742158), and the Wellcome Trust (grant 206194).

The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Copyright: © 2019 Lintusaari J *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

How to cite this article: Lintusaari J, Blomstedt P, Rose B *et al.* **Resolving outbreak dynamics using approximate Bayesian computation for stochastic birth–death models [version 2; peer review: 1 approved, 1 approved with reservations]** Wellcome Open Research 2019, 4:14 (<https://doi.org/10.12688/wellcomeopenres.15048.2>)

First published: 25 Jan 2019, 4:14 (<https://doi.org/10.12688/wellcomeopenres.15048.1>)

two BD processes and one pure birth process that have epidemiologically based interpretations. As in a standard BD process, these events are assumed to be independent of one another and to occur at specific rates. The time between two events is assumed to follow the exponential distribution specified by the rate of occurrence, causing the number of events to follow the Poisson distribution. The timescale considered here is one calendar year. The evolution of the infectious population is simulated by drawing events according to their rates.

Building upon the BD process, the simulated population carries auxiliary information. At birth, each case is assigned a cluster index that represents the specific genetic fingerprint of the pathogen and determines the cluster the case belongs to. The simulated output includes the cluster indexes that are recorded when cases are observed.

We will now explain our model in more detail and point out differences between it and the model of Tanaka *et al.*².

First, we assume that observations are collected within a given time interval that matches that of the observed data. In the case of the San Francisco Bay data, the length of this interval is two years¹¹. Observations are collected from the simulated process after a sufficient warmup period so that the process can be expected to have reached stable properties. This procedure is visualized in Figure 1.

In the figure, the dashed lines are the balance values. The population sizes fluctuate around them after the process has matured. Both populations surpass their balance values at least once by the 22-year mark. The observation period is the green patch. The grey line shows the number of observations collected during each year of the simulation. The number of observations from the observation period and the clustering structure of the observations are used in the inference of the epidemiological

parameters. A patient becomes observed in the study with probability p_{obs} . Our model makes the simplifying assumption that both being observed and ceasing to be infectious are combined under the death event in the simulation. This is based on the assumption that a typical patient is treated promptly after being diagnosed¹³, but we still allow for the possibility that some patients do not comply with treatment and remain infectious (see below). In contrast to the model of Tanaka *et al.*², there is no separate observation sampling phase, nor is there a prior estimate for the underlying population size.

We introduce a burden parameter β that reflects the rate at which new active TB cases with a previously unseen pathogen fingerprint appear in the community. This is the pure birth process of the model, and it represents reactivation of TB from latent cases as well as new pathogen fingerprints introduced by immigration. In the simulation, each such case receives a new cluster index that has not been assigned to any earlier case. Unlike Tanaka *et al.*², we do not explicitly model mutations. Instead, we assume they occur during the latent phase of infection over the years¹¹. This decision was partially motivated by the fact that Aandahl *et al.*¹⁰ found the mutation rate parameter from 2 to be non-identifiable from the fingerprint data, and they consequently fixed that value to a constant.

We introduce two distinct birth–death processes for cases that are either *compliant* or *non-compliant* with treatment. These birth–death processes are parametrized with birth rates τ_i and death rates δ_i , where $i = 1$ denotes the non-compliant population and $i = 2$ the compliant population. A significant number of cases in the largest clusters observed by Small *et al.*¹¹ corresponded to non-compliant patients who stayed infectious for several months and belonged to subgroups under increased risk of rapid development of active TB due to conditions such as AIDS and substance abuse. Patients who are compliant with therapy typically cease being infectious

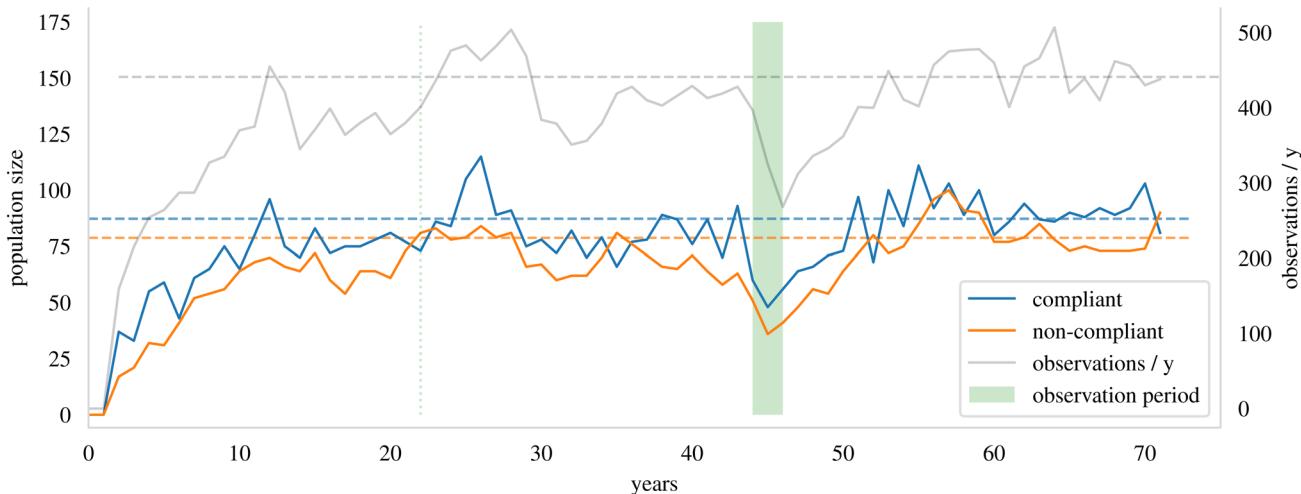


Figure 1. An illustration of simulated compliant and non-compliant populations as observed at the end of each year. Note that sampling can be done at any point once the model has stabilized; the drop in population sizes at the sampling point in this figure is purely coincidental.

$$\begin{aligned}
R_1 &\sim \text{Unif}(1.01, 20), \\
R_2 | R_1 &\sim \text{Unif}(0.01, (1 - 0.05 \cdot R_1)/0.95), \\
\text{and} \\
t_1 &\sim \text{Unif}(0.01, 30),
\end{aligned} \tag{7}$$

Given the observed data, we set the following additional constraints to optimize computation:

$$\begin{aligned}
\hat{n}_{obs} < 350, \\
\text{and} \\
\tau_1 < 40.
\end{aligned} \tag{8}$$

We verified that these constraints have a negligible effect on the acquired estimates. Their function is to prevent simulations with extremely unlikely parameter values, which saves a considerable amount of computation time. As a result of these constraints, all obtained estimates of R_1 are smaller than 15. [Figure 2](#) shows the samples drawn from the priors under these conditions.

3.1.2 Summary statistics. The summary statistics used in earlier approaches (e.g. [2](#) and [12](#)) are not directly applicable to our model. This is due to differences between the models that cause, for example, the number of observations in the sample to vary rather than being fixed. However, the previous studies' summaries did prove to be a good starting point for the

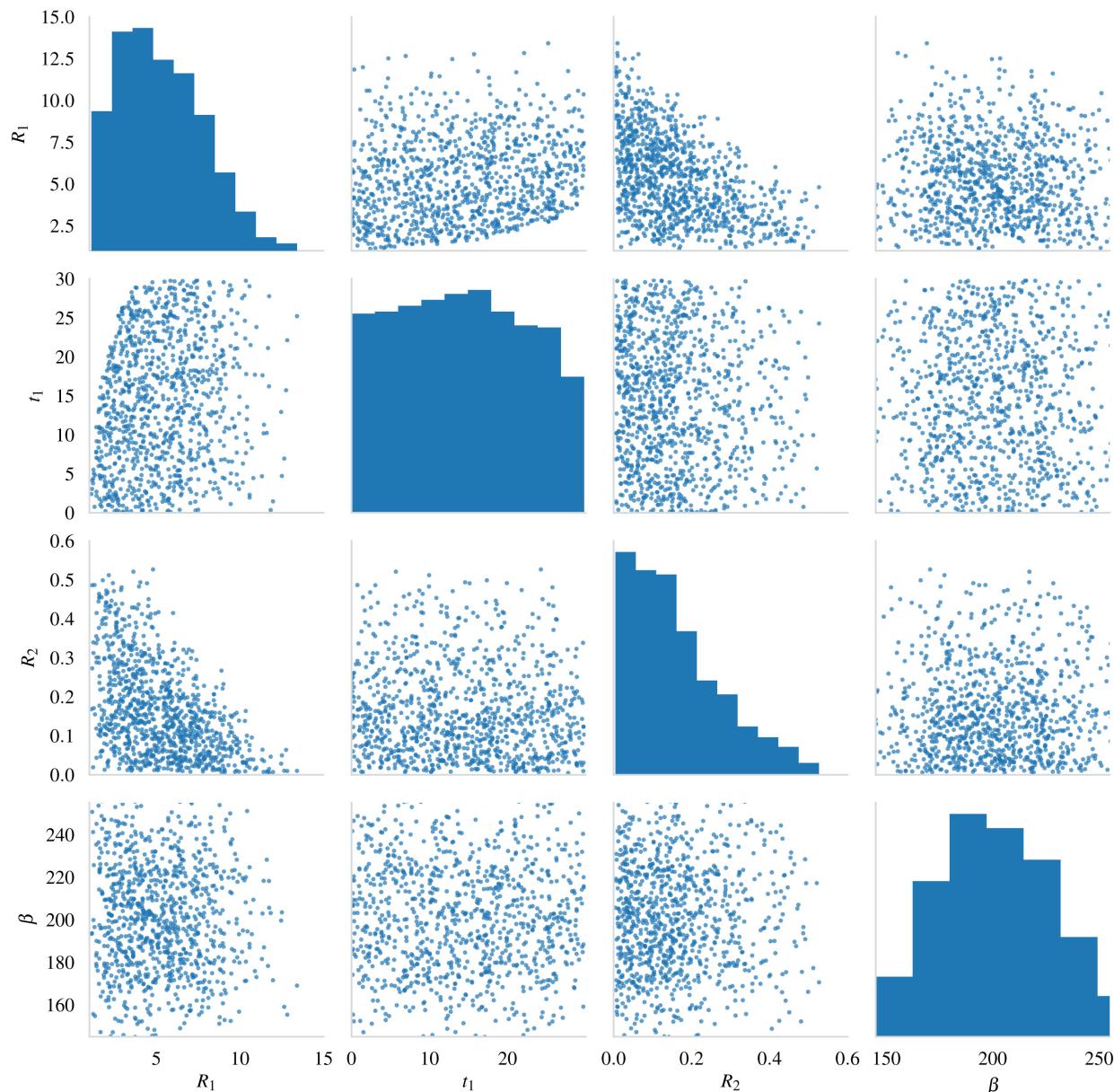


Figure 2. A scatter matrix of samples from the prior.

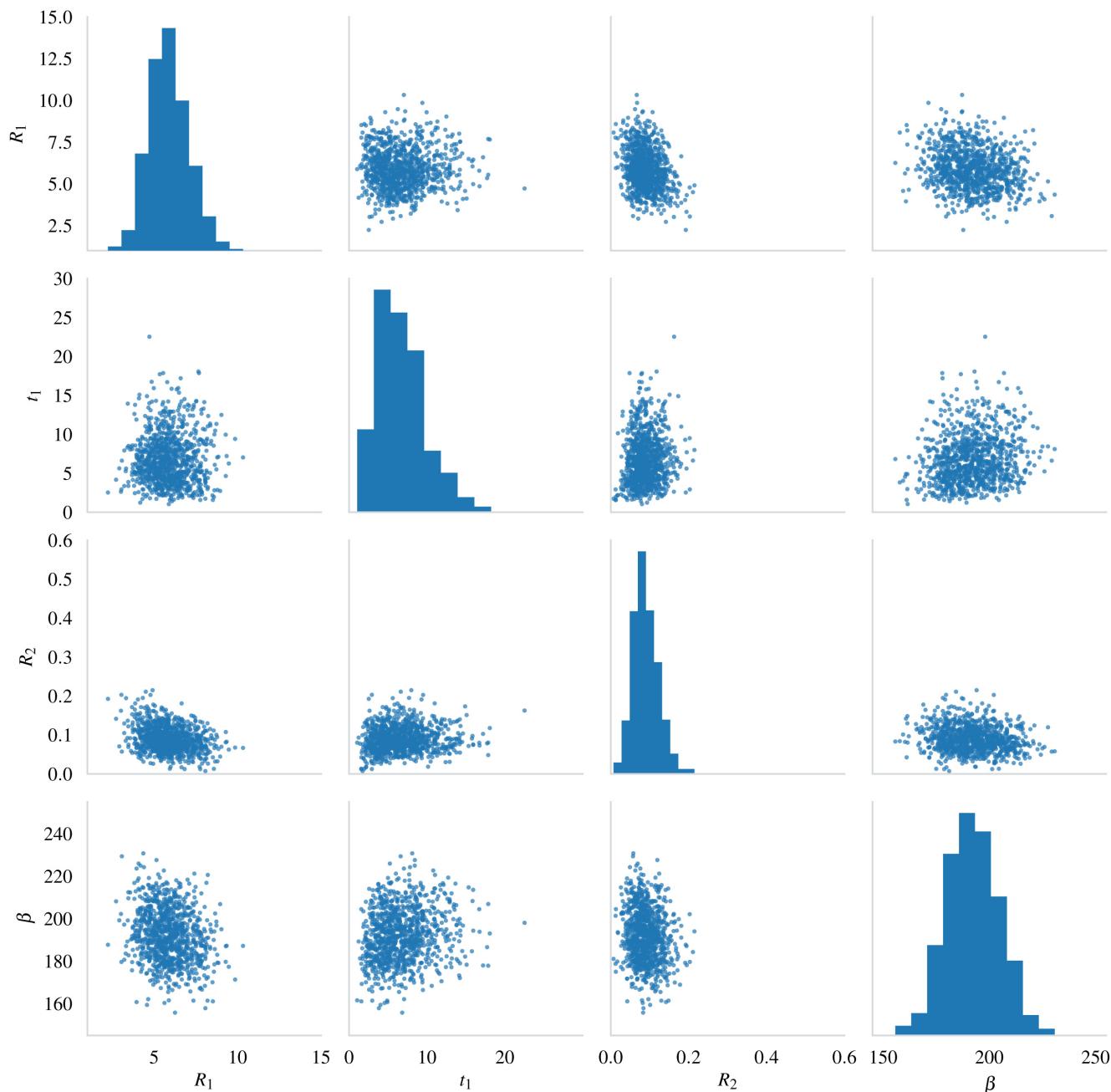


Figure 3. Posterior sample of size 1000 from the approximate posterior distribution $\tilde{p}(R_1, t_1, R_2, \beta | y_0)$ plotted as a scatter matrix. Compare to the prior in Figure 2.

Table 2. Posterior summaries.

Parameter	Mean	Median	95% CI
R_1	5.88	5.79	(3.68, 8.16)
t_1	6.74	6.25	(1.57, 12.9)
R_2	0.09	0.09	(0.03, 0.15)
β	192	192	(170, 216)

Table 3 lists the MAE and MdAE with the 95% error upper percentile for each parameter estimate. This information is useful for quantifying how much each estimate deviates from the actual parameter value on average. The burden rate (β) and the reproductive number of the non-compliant population (R_1) have the smallest relative MAEs: 4.0% and 14.9%, respectively. The reproductive number of the compliant population (R_2) and the net transmission rate of the non-compliant population (t_1) have MAEs of 29.5% and 44.2%, respectively.

The MAE of the latter seems rather high. The 95% percentile indicates that in 5% of the trials, the error was substantial. Further investigation of this issue shows that for some of the synthetic datasets, t_1 is not identifiable, meaning that the synthetic data in those cases is not informative enough to produce a clear mode for the parameter. R_2 suffered slightly from the same problem. This kind of situation, where some of the synthetic datasets turn out uninformative, is rather common when little data is available. Because of these exceptions, the MdAE might be a more appropriate measure than the MAE, as the former is not as heavily influenced by the results of non-identifiable datasets in trials. The relative MdAE errors for R_2 and t_1 were 21.9% and 32.1%, respectively.

Figure 4 visualizes the estimated vs. actual values of each of the parameters.

Though t_1 is only weakly identifiable, the results of our simulations indicate that the set of epidemiological parameters we

have analyzed is identifiable for the San Francisco Bay dataset. Our simulations suggest that a structural model issue could be at fault for the weak identifiability of t_1 , as this can arise as a consequence of the generating stochastic process producing a relatively flat cluster distribution. Fortunately, β , R_1 and R_2 , all of which provide more valuable epidemiological insight than t_1 , are robust against this identifiability issue.

The coverage property¹⁷ is used to assess the reliability of the inference by checking whether the spreads of the acquired posterior distributions are accurate. Given a critical level α , the true parameter value should be outside the $1-\alpha$ credible interval of the posterior with probability α . We carried out our coverage analysis as follows.

First, we used rejection sampling to produce a sample for the posterior from the observed data. From this posterior, we sampled 1000 parameter vectors (with replacement) for the trials. For each of these 1000 vectors, we simulated synthetic

Table 3. Mean and median absolute errors for 1000 trials with synthetic data from the posterior.

Parameter	MAE	Relative MAE ²	MdAE	Relative MdAE	95% percentile
R_1	0.85	14.9%	0.72	12.6%	2.00
t_1	2.68	44.2%	1.98	32.1%	7.66
R_2	0.024	29.5%	0.018	21.9%	0.07
β	7.6	4.0 %	6.1	3.1%	19.8

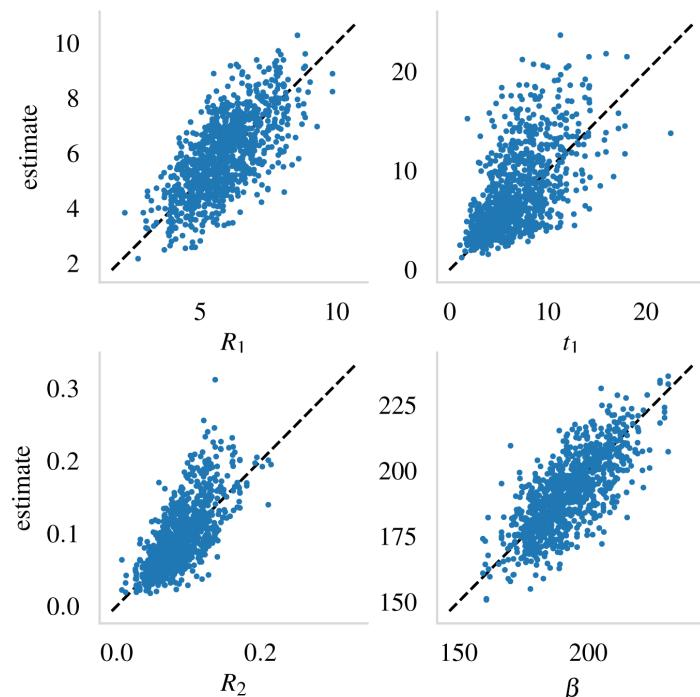


Figure 4. The estimates from the 1000 trials plotted against their true values. The black dashed line shows the 1:1 correspondence.

16. Nunes MA, Balding DJ: **On optimal selection of summary statistics for approximate Bayesian computation.** *Stat Appl Genet Mol Biol.* 2010; **9**(1): Article34.
[PubMed Abstract](#) | [Publisher Full Text](#)
17. Wegmann D, Leuenberger C, Excoffier L: **Efficient approximate Bayesian computation coupled with Markov chain Monte Carlo without likelihood.** *Genetics.* 2009; **182**(4): 1207–18.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
18. Lintusaari J, Gutmann MU, Dutta R, et al.: **Fundamentals and Recent Developments in Approximate Bayesian Computation.** *Syst Biol.* 2017; **66**(1): e66–e82.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Open Peer Review

Current Peer Review Status:  

Version 2

Reviewer Report 07 October 2019

<https://doi.org/10.21956/wellcomeopenres.16856.r36360>

© 2019 Gascuel O. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



Olivier Gascuel

Unité Bioinformatique Evolutive, C3BI UMR 3756 Institut Pasteur & CNRS, Paris, France

The authors have made changes that further improved the clarity and the flow of the manuscript and addressed our comments and suggestions. We propose final approval of the paper. We would like to suggest several minor changes not requiring any further analysis.

1. While adding a workflow of the study helps greatly, we think that adding a flow diagram of the model (*i.e.* the flow of individuals between different states while pointing out individual rates, e.g. see Figure 1 at <https://institutefordiseasemodeling.github.io/Documentation/general/model-seir.html> for SEIR model) would help the readers to grasp quickly the mathematical model.
2. We believe that adding some parts of the original, more extensive, description of the Figure 1 explaining the meaning of dashed lines (balance values) should be kept.
3. "It should be noted that the summaries chosen here do not consider global sufficiency. In cases where the dataset is very different from the San Francisco data, a modified set of summaries should probably be considered."

With respect to this remark, we believe an extensive list of the summary statistics that you tried before identifying the final set, would help further research and adaptation of your method to similar studies. Such a list might be added as Supplementary Material.

Hope this helps, sincerely,
Jakub Voznica, Anna Zhukova and Olivier Gascuel

Competing Interests: No competing interests were disclosed.

I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.

Version 1

Reviewer Report 12 April 2019

<https://doi.org/10.21956/wellcomeopenres.16417.r34750>

© 2019 Beaumont M. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**Mark Beaumont**

School of Biological Sciences, University of Bristol, Bristol, UK

This interesting paper covers an area - TB epidemiology - that has been the subject of a number of papers that have used approximate Bayesian computation (ABC). These studies have also involved some full-likelihood MCMC solutions for the same models. The present paper carries out an ABC analysis with an alternative modelling framework, which provides some satisfactory solutions to some problems that had been previously noted.

I generally have few quibbles with the ABC analysis. My only main query with the paper, in my ignorance of TB epidemiology, is to what extent the new modelling framework gives a reasonable representation of the underlying biology. I note that the original Tanaka *et al.* paper explicitly emphasised the importance of modelling the mutational structure of clusters (3rd paragraph of the introduction). By contrast the present paper gives no justification for dropping the mutational modelling. My reason for querying this is that, as I understand it, Tanaka had a previous history of detailed work on TB epidemiology, prior to the 2006 paper, including co-authorship with Peter Small, and so presumably was able to put in the benefit of that experience into the paper. Therefore I recommend that this aspect be much better justified and discussed.

On the ABC side, my main query is to what extent the authors are confident about identifiability of their parameters. Particularly, since they seem to suggest one of their parameters is not identifiable (discussed more in specific points below). In a model-free setting, identifiability is demonstrated through simulation, rather than analytically. Obviously, if non-identifiability is shown this naturally leads to some questions about the summary statistics etc. as well as the structure of the model itself. But, with informative priors, some parameters that are only jointly identifiable can appear to be identifiable marginally - in a population genetics context the apparent identifiability of N and $\lambda\mu$ with informative priors is a case in point, when only their product is identifiable. Again, this needs a bit more discussion than in the present paper.

Specific Comments:

- Introduction, first paragraph: "genotype fingerprints". Some more discussion of this would be useful with regard to my point above. Presumably what concerned Tanaka *et al.* is that multiple outbreaks can involve the same cluster, and that different clusters (due to mutation?) could arise from the same outbreak.
- Model, 4th paragraph: " $p_{\{obs\}}$ " - does the assumption of being observed lead to ceasing to be infectious fit with the compliant/non-compliant distinction, two paragraphs further down?

- Model, 5th paragraph: Note my main query.
- Summary statistics, paragraph 2: It might be helpful to emphasise that a 'cluster' here is assumed to be a new active TB case. Presumably many of these summary statistics are highly correlated with the parameter β ?
- Figure 2/3: I wonder whether these might be better in the supp. text, and replaced with a single figure with HPD contours for the prior and posterior.
- Summary statistics, last paragraph: "It is good to note". Do the authors mean that? Or rather do they mean "It should be noted"? Presumably it is not good to be not sufficient. More generally, is there an argument to use projections as in Fearnhead and Prangle (2012¹), which are generally straightforward to apply? I think there are good reasons (Fearnhead and Prangle, 2012¹; Li and Fearnhead, 2018²) for expecting the optimal number of summary statistics to be the same as the number of parameters, thus reducing the effect of the 'curse of dimensionality'.
- Results, 3rd paragraph: coverage property. The analysis seems fine, but the authors skirt some details, worth noting. They use 'true' values from the ABC posterior. The Wegmann *et al.* paper, following Cook *et al.*, simulated 'true' values from the prior, for which coverage is indeed uniform. It is not so obvious that coverage from the ABC posterior should also be uniform, but this is demonstrated (I think for the first time) in Prangle *et al.* (2014³) (at least for any interval in the prior predictive distribution of summary statistics, including the interval from which the ABC posterior is computed).
- Results, 4th paragraph: non-identifiability of t_1 . This observation seems at variance with what is stated in the abstract. Is this a summary statistic issue? Or a structural model issue?
- Figure 5: These results look convincing. Note that because of the need for a tolerance interval ABC coverage is not expected to be perfect (Fearnhead and Prangle, 2012¹).

References

1. Fearnhead P, Prangle D: Constructing summary statistics for approximate Bayesian computation: semi-automatic approximate Bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2012; **74** (3): 419-474 [Publisher Full Text](#)
2. Li W, Fearnhead P: On the asymptotic efficiency of approximate Bayesian computation estimators. *Biometrika*. 2018; **105** (2): 285-299 [Publisher Full Text](#)
3. Prangle D, Blum M, Popovic G, Sisson S: Diagnostic tools for approximate Bayesian computation using the coverage property. *Australian & New Zealand Journal of Statistics*. 2014; **56** (4): 309-329 [Publisher Full Text](#)

Is the rationale for developing the new method (or application) clearly explained?

Partly

Is the description of the method technically sound?

Yes

Are sufficient details provided to allow replication of the method development and its use by others?

Yes

If any results are presented, are all the source data underlying the results available to ensure full reproducibility?

Yes

Are the conclusions about the method and its performance adequately supported by the findings presented in the article?

Partly

Competing Interests: No competing interests were disclosed.

Reviewer Expertise: Statistical population genetics

I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.

Author Response 30 Jul 2019

Jarno Lintusaari, Aalto University, Espoo, Finland

We thank the reviewers for their useful comments that allowed us to improve the paper. Below we provide detailed responses to the issues brought up. Our responses are written in italics.

This interesting paper covers an area - TB epidemiology - that has been the subject of a number of papers that have used approximate Bayesian computation (ABC). These studies have also involved some full-likelihood MCMC solutions for the same models. The present paper carries out an ABC analysis with an alternative modelling framework, which provides some satisfactory solutions to some problems that had been previously noted.

I generally have few quibbles with the ABC analysis. My only main query with the paper, in my ignorance of TB epidemiology, is to what extent the new modelling framework gives a reasonable representation of the underlying biology. I note that the original Tanaka et al. paper explicitly emphasised the importance of modelling the mutational structure of clusters (3rd paragraph of the introduction). By contrast the present paper gives no justification for dropping the mutational modelling. My reason for querying this is that, as I understand it, Tanaka had a previous history of detailed work on TB epidemiology, prior to the 2006 paper, including co-authorship with Peter Small, and so presumably was able to put in the benefit of that experience into the paper. Therefore I recommend that this aspect be much better justified and discussed.

This is a very relevant point. However, the subsequent joint work by the Tanaka et al. authors and Tanya Stadler notified that the mutation parameter appears non-identifiable from the fingerprint data and used a fixed value obtained from the literature in their most recent paper. We have now included a discussion about this and the rationale for excluding an explicit mutation rate parameter in our model.

On the ABC side, my main query is to what extent the authors are confident about identifiability of their parameters. Particularly, since they seem to suggest one of their parameters is not identifiable (discussed more in specific points below). In a model-free setting, identifiability is demonstrated through simulation, rather than analytically. Obviously, if non-identifiability is shown this naturally leads to some questions about the summary statistics etc. as well as the structure of the model itself. But, with informative priors, some parameters that are only jointly identifiable can appear to be identifiable marginally - in a population genetics context the apparent identifiability of N and μ with informative priors is a case in point, when only their product is identifiable. Again, this needs a bit more discussion than in the present paper.

We agree that an additional discussion is in place and have added such in the revision. Our simulation experiments reported in the paper suggest that the key epidemiological parameters are indeed identifiable for the SF Bay data set, even if the net transmission rate may remain only weakly identifiable.

Specific Comments:

Introduction, first paragraph: "genotype fingerprints". Some more discussion of this would be useful with regard to my point above. Presumably what concerned Tanaka et al. is that multiple outbreaks can involve the same cluster, and that different clusters (due to mutation?) could arise from the same outbreak.

We have added further discussion. Noting the slow mutation rate of TB, it is highly unlikely that multiple clusters would arise from the same outbreak within a relatively short timespan.

Model, 4th paragraph: " p_{obs} " - does the assumption of being observed lead to ceasing to be infectious fit with the compliant/non-compliant distinction, two paragraphs further down?

Good point, there was a sloppy phrasing in the 4th paragraph. We have now revised the text to be in line with the later paragraph.

Model, 5th paragraph: Note my main query.

As noted in our response to the main query item, we have now edited the text accordingly.

Summary statistics, paragraph 2: It might be helpful to emphasise that a 'cluster' here is assumed to be a new active TB case. Presumably many of these summary statistics are highly correlated with the parameter β ?

Excellent remarks, we have added further clarification about this.

Figure 2/3: I wonder whether these might be better in the supp. text, and replaced with a single figure with HPD contours for the prior and posterior.

We do appreciate this suggestion, however, as noted in the response to R1, the first author who had the main responsibility for all aspects of the presented work has already graduated and left academia, as has the second author, so neither of the two are able to contribute to further work related to this paper. We would thus prefer keeping the two figures as in their current versions.

Summary statistics, last paragraph: "It is good to note". Do the authors mean that? Or rather do they mean "It should be noted"? Presumably it is not good to be not sufficient.

The reviewer has a correct interpretation, this was a typo and is now fixed.

More generally, is there an argument to use projections as in Fearnhead and Prangle (2012¹), which are generally straightforward to apply? I think there are good reasons (Fearnhead and Prangle, 2012¹; Li and Fearnhead, 2018²) for expecting the optimal number of summary statistics to be the same as the number of parameters, thus reducing the effect of the 'curse of dimensionality'.

It is indeed correct that the number of summary statistics is generally expected to match the dimensionality of the parameter space. As noted in the response to R1, the summary statistics were iteratively defined by trialing inference on synthetic data from the model. The final set of statistics was settled on after extensive test simulations showing appropriate behavior. Note also that we had previous experience about the behavior of various summary statistics from the earlier Lintusaari et al. Genetics 2016 article examining inference for a different model but the same data.

Results, 3rd paragraph: coverage property. The analysis seems fine, but the authors skirt some details, worth noting.

Detailed description of the simulations used to study the coverage property has been added.

They use 'true' values from the ABC posterior. The Wegmann et al. paper, following Cook et al., simulated 'true' values from the prior, for which coverage is indeed uniform. It is not so obvious that coverage from the ABC posterior should also be uniform, but this is demonstrated (I think for the first time) in Prangle et al. (2014³) (at least for any interval in the prior predictive distribution of summary statistics, including the interval from which the ABC posterior is computed).

Good point, we have added reference to Prangle et al. in the relevant part of the text.

Results, 4th paragraph: non-identifiability of t_1. This observation seems at variance with what is stated in the abstract. Is this a summary statistic issue? Or a structural model issue?

We have edited the text to clarify the potential weak identifiability of t_1. Our simulations suggest that it is a structural model issue such that the parameter sometimes becomes only weakly identifiable when the generating stochastic process happens to result in a particularly flat distribution of clusters. Encouragingly, the other parameters, which are the most relevant ones from the epidemiological perspective, do appear fairly robustly identifiable even if t_1 would have a flat posterior.

Figure 5: These results look convincing. Note that because of the need for a tolerance interval ABC coverage is not expected to be perfect (Fearnhead and Prangle, 2012¹).

Good point, we have added a remark about this to the revised text.

Competing Interests: No competing interests were disclosed.

Reviewer Report 13 February 2019

<https://doi.org/10.21956/wellcomeopenres.16417.r34679>

© 2019 Gascuel O et al. This is an open access peer review report distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.



Jakub Voznica

Unité Bioinformatique Evolutive, C3BI UMR 3756 Institut Pasteur & CNRS, Paris, France

Anna Zhukova

Unité Bioinformatique Evolutive, C3BI UMR 3756 Institut Pasteur & CNRS, Paris, France

Olivier Gascuel

Unité Bioinformatique Evolutive, C3BI UMR 3756 Institut Pasteur & CNRS, Paris, France

Article summary

The article describes a new model of TB outbreak in San Francisco Bay area that overcomes the non-identifiability/dependency on the assumed population size of the reproductive number R in the generic birth-death-mutation model by Tanaka *et al.* The new model considers two compartments, for compliant and non-compliant subpopulations, and combines two birth-death processes (for each of the compartments) with a pure-birth process that creates new TB transmission clusters (i.e. a new individual with a new RFLP pattern that is further transmitted). This pure-birth process replaces mutation in Tanaka's model and corresponds to migration or reactivating of a latent TB. The rate corresponding to the pure-birth process is referred as the burden rate. At each (non-burden) birth event (i.e. TB transmission) the compartment of the newly infected individual is assigned to non-compliant or compliant with the probability p_1 or $(1 - p_1)$ correspondingly. At each death (i.e. becoming non-infectious) event the individual is sampled with the probability p_{obs} .

Overall, the proposed model has 7 parameters: the burden rate, 2 birth rates, 2 death rates, and 2 probabilities (p_1 and p_{obs}). However, 3 of them (compliant death rate, p_{obs} and p_1) were fixed based on the estimates from the literature, therefore leaving 4 parameters to be estimated, expressed in terms of two reproductive numbers, i.e. birth to death rate ratios for the corresponding compartments, the non-compliant net transmission rate (difference between the birth and the death rates), and the burden rate. Priors and additional constraints on the rates were set to avoid biological meaningless of the simulations.

The simulator was implemented for the proposed model and parameter estimation was performed for the data collected in SF Bay area in 1991-92 (Small *et al.*) with ABC, based on 1000 parameter values sampled with rejection from 6M simulations, using 8 (weighted) summary statistics:

1. the number of observations
2. the total number of clusters
3. the relative number of singleton clusters
4. the relative number of clusters of size two
5. the size of the largest cluster
6. the mean of the successive difference in size among the four largest clusters
7. the number of months from the first observation to the last
8. the number of months when at least one observation was made.

The new model not only allowed for estimation of the aforementioned parameters (posterioris are well concentrated within but far from the edges of the priors) but also of the balance subpopulation sizes (at the equilibrium state when infected subpopulations neither shrink nor grow). The estimates differ from those done with the birth-death-mutation model, and are potentially better aligned with the epidemiological knowledge on TB in the area.

The coverage property (accuracy of the spread of the acquired posterior) of the estimator was further tested on 1000 parameter values drawn from the posterior, giving satisfactory results for the critical level of .05 (the true parameter values were outside of the .95 credible interval of the posterior with probability less than .05).

General comments

The article reads well, the model, rationale behind it, its assumptions and advantages over the previous TB model are explained in a clear and convincing way. It is a valuable addition to TB research, and we believe that the article should be accepted.

Having little knowledge on TB (but on ABC), we feel like the article could benefit from a more detailed discussion of the obtained estimates. For example, is there any literature/other data supporting the estimated subpopulation sizes?

We also point out a few technicalities that could be explained in more detail (see below).

Technical comments

A flow diagram of the model could facilitate the model understanding for the reader.

Additional sensitivity analysis of the model while varying pre-fixed parameter values (of compliant death rate, p-obs and p1) might add confidence in author's findings.

Page 4: "*The observations are collected from the simulated process after a sufficient warm-up period, so that the process can be expected to have reached stable properties (exemplified in Figure 1).*"

In Figure 1 the warm-up seems to be achieved already after 15 years, however the observation period is chosen around 45 years, where there is a drop of population sizes. Is it a coincidence? How is the start of the observation period selected?

Page 5: "*We used the Engine for Likelihood-Free Inference (ELFI)...*"

The authors might detail what kind of inference was used: Is it a pure distance/rejection-based approach? Or do you use some regression tool, random forest, LASSO, neural network or other? How was the technique selected?

Page 5: "*Based on the details in Small et al. describing the San Francisco Bay area TB data, there were 585 confirmed cases of TB of which 487 were included in the study. To account for the cases that were not included in the study, we fix the probability of becoming observed to p-obs = 0.8*"

If we understand correctly the p-obs is calculated as 487/585, but what about potentially unknown cases of TB in the SF Bay area? Is it assumed that all the existing TB cases are known?

Page 5: It is not very clear why these particular summary statistics were selected, e.g. "*the mean of the*

successive difference in size among the four largest clusters”

Why not 3 or 5, etc.? Were for example other statistics tested, which performed worse?

The name of the last statistic (“*the number of months when at least one observation was made*”) is rather confusing. In table 1 it has a slightly different name: “*the number of months that at least one observation was made from the largest cluster*”. Does it mean *the time when the first observation from the largest cluster was made*?

Page 7: “*The chosen summary statistics and weights were found to perform well in the evaluation of the model in Subsection .*”

The subsection number is missing.

Page 7: “*The resulting threshold for the acquired sample was $\epsilon = 31.7$ with the smallest distance being 12.5.*”

How were the threshold and distance values selected?

Is the rationale for developing the new method (or application) clearly explained?

Yes

Is the description of the method technically sound?

Yes

Are sufficient details provided to allow replication of the method development and its use by others?

Yes

If any results are presented, are all the source data underlying the results available to ensure full reproducibility?

Yes

Are the conclusions about the method and its performance adequately supported by the findings presented in the article?

Yes

Competing Interests: No competing interests were disclosed.

We confirm that we have read this submission and believe that we have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however we have significant reservations, as outlined above.

Author Response 30 Jul 2019

Jarno Lintusaari, Aalto University, Espoo, Finland

We thank the reviewers for their useful comments that allowed us to improve the paper. Below we provide detailed responses to the issues brought up. Our responses are written in italics.

R1:

The article describes a new model of TB outbreak in San Francisco Bay area that overcomes the

non-identifiability/dependency on the assumed population size of the reproductive number R in the generic birth-death-mutation model by Tanaka et al. The new model considers two compartments, for compliant and non-compliant subpopulations, and combines two birth-death processes (for each of the compartments) with a pure-birth process that creates new TB transmission clusters (i.e. a new individual with a new RFLP pattern that is further transmitted). This pure-birth process replaces mutation in Tanaka's model and corresponds to migration or reactivating of a latent TB. The rate corresponding to the pure-birth process is referred as the burden rate. At each (non-burden) birth event (i.e. TB transmission) the compartment of the newly infected individual is assigned to non-compliant or compliant with the probability p1 or (1 - p1) correspondingly. At each death (i.e. becoming non-infectious) event the individual is sampled with the probability p-obs.

Overall, the proposed model has 7 parameters: the burden rate, 2 birth rates, 2 death rates, and 2 probabilities (p1 and p-obs). However, 3 of them (compliant death rate, p-obs and p1) were fixed based on the estimates from the literature, therefore leaving 4 parameters to be estimated, expressed in terms of two reproductive numbers, i.e. birth to death rate ratios for the corresponding compartments, the non-compliant net transmission rate (difference between the birth and the death rates), and the burden rate. Priors and additional constraints on the rates were set to avoid biological meaningless of the simulations.

The simulator was implemented for the proposed model and parameter estimation was performed for the data collected in SF Bay area in 1991-92 (Small et al.) with ABC, based on 1000 parameter values sampled with rejection from 6M simulations, using 8 (weighted) summary statistics:

1. the number of observations
2. the total number of clusters
3. the relative number of singleton clusters
4. the relative number of clusters of size two
5. the size of the largest cluster
6. the mean of the successive difference in size among the four largest clusters
7. the number of months from the first observation to the last
8. the number of months when at least one observation was made.

The new model not only allowed for estimation of the aforementioned parameters (posterioris are well concentrated within but far from the edges of the priors) but also of the balance subpopulation sizes (at the equilibrium state when infected subpopulations neither shrink nor grow). The estimates differ from those done with the birth-death-mutation model, and are potentially better aligned with the epidemiological knowledge on TB in the area.

The coverage property (accuracy of the spread of the acquired posterior) of the estimator was further tested on 1000 parameter values drawn from the posterior, giving satisfactory results for the critical level of .05 (the true parameter values were outside of the .95 credible interval of the posterior with probability less than .05).

General comments

The article reads well, the model, rationale behind it, its assumptions and advantages over the previous TB model are explained in a clear and convincing way. It is a valuable addition to TB research, and we believe that the article should be accepted.

We thank the reviewers for their highly positive comments about our work.

Having little knowledge on TB (but on ABC), we feel like the article could benefit from a more detailed discussion of the obtained estimates. For example, is there any literature/other data supporting the estimated subpopulation sizes?

The estimates are well aligned with the epidemiological discussion in the original NEJM paper introducing the fingerprint data. We now point this out more carefully in the revised version.

We also point out a few technicalities that could be explained in more detail (see below).

Technical comments

A flow diagram of the model could facilitate the model understanding for the reader.

A flow diagram has been added as a supplementary figure to accompany the final version.

Additional sensitivity analysis of the model while varying pre-fixed parameter values (of compliant death rate, p-obs and p1) might add confidence in author's findings.

We feel that the current sensitivity analysis is quite sufficient and fulfils its purpose to demonstrate stability of the estimates for data akin the San Francisco Bay observations. The first author who had the main responsibility for all aspects of the presented work has already graduated and left academia, so he is not able to contribute to further work related to this paper which limits our possibilities for performing extensive additional simulations.

Page 4: "The observations are collected from the simulated process after a sufficient warm-up period, so that the process can be expected to have reached stable properties (exemplified in Figure 1)."

In Figure 1 the warm-up seems to be achieved already after 15 years, however the observation period is chosen around 45 years, where there is a drop of population sizes. Is it a coincidence? How is the start of the observation period selected?

Figure 1 is intended only as a schematic numerical example to assist the reader in understanding the underlying logic of the model and the sampling assumptions. The drop is thus coincidental. This is now properly noted in the revision.

Page 5: "We used the Engine for Likelihood-Free Inference (ELFI)..."

The authors might detail what kind of inference was used: Is it a pure distance/rejection-based approach? Or do you use some regression tool, random forest, LASSO, neural network or other? How was the technique selected?

As stated in the paper, we used pure rejection sampling (we sampled 1000 parameter values with rejection sampling from a total of 6M simulations). The main reasons for choosing this basic ABC approach were: 1) we implemented a computationally efficient vectorized Python version of the

simulator which facilitated the use of a large number of simulations, 2) at the time the project was initiated, ELFI had yet no implementation of the Bayesian optimization procedure for non-uniform priors. Such a prior was essential for the model structure and straightforward to consider in a pure ABC rejection sampler, hence the choice for inference method was well motivated. These reasons are now more clearly stated in the revision.

Page 5: “Based on the details in Small et al. describing the San Francisco Bay area TB data, there were 585 confirmed cases of TB of which 487 were included in the study. To account for the cases that were not included in the study, we fix the probability of becoming observed to p-obs = 0.8” If we understand correctly the p-obs is calculated as 487/585, but what about potentially unknown cases of TB in the SF Bay area? Is it assumed that all the existing TB cases are known?

For epidemiological reasons it is unlikely that any substantial numbers of active TB cases were unknown to the public health officials, hence it is unlikely that these would have a non-negligible contribution to the observed outbreaks. Given the severity of TB and the protocols followed by public health officials most active cases are expected to have been traced. We have now stated this more explicitly in the revision.

Page 5: It is not very clear why these particular summary statistics were selected, e.g. “the mean of the successive difference in size among the four largest clusters”

Why not 3 or 5, etc.? Were for example other statistics tested, which performed worse? The name of the last statistic (“the number of months when at least one observation was made”) is rather confusing. In table 1 it has a slightly different name: “the number of months that at least one observation was made from the largest cluster”. Does it mean the time when the first observation from the largest cluster was made?

We have edited the text to make the summary statistic definitions unambiguous. The summary statistics were iteratively defined by trialing inference on synthetic data from the model. The final set of statistics was settled on after extensive test simulations showing appropriate behavior. Note also that we had previous experience about the behavior of various summary statistics from the earlier Lintusaari et al. Genetics 2016 article examining inference for a different model but the same data.

Page 7: “The chosen summary statistics and weights were found to perform well in the evaluation of the model in Subsection .”

The subsection number is missing.

The subsection number was missing due to the submission template and will be visible in the final typeset version.

Page 7: “The resulting threshold for the acquired sample was $\epsilon = 31.7$ with the smallest distance being 12.5.”

How were the threshold and distance values selected?

As for the summary statistics, the threshold was settled by extensive trialing of inference on

synthetic data from the model to identify a threshold striking a good balance between runtimes and acceptance rate and the resulting Monte Carlo error rate. This is now more appropriately reported in the revision.

Competing Interests: No competing interests were disclosed.