

PUBLISHED VERSION

Black, Andrew James; Ross, Joshua Vincent

Estimating a Markovian epidemic model using household serial interval data from the early phase of an epidemic, *PLoS One*, 2013; 8(8):e73420.

© 2013 Black, Ross. 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 author and source are credited.

PERMISSIONS

<http://www.plosone.org/static/license>

Open-Access License



No Permission Required

PLOS applies the [Creative Commons Attribution License](#) (CCAL) to all works we publish (read the [human-readable summary](#) or the [full license legal code](#)). Under the CCAL, authors retain ownership of the copyright for their article, but authors allow anyone to download, reuse, reprint, modify, distribute, and/or copy articles in PLOS journals, so long as the original authors and source are cited. **No permission is required from the authors or the publishers.**

In most cases, appropriate attribution can be provided by simply citing the original article (e.g., Kaltenbach LS et al. (2007) Huntingtin Interacting Proteins Are Genetic Modifiers of Neurodegeneration. *PLOS Genet* 3(5): e82. doi:10.1371/journal.pgen.0030082). If the item you plan to reuse is not part of a published article (e.g., a featured issue image), then please indicate the originator of the work, and the volume, issue, and date of the journal in which the item appeared. For any reuse or redistribution of a work, you must also make clear the license terms under which the work was published.

This broad license was developed to facilitate open access to, and free use of, original works of all types. Applying this standard license to your own work will ensure your right to make your work freely and openly available. Learn more about [open access](#). For queries about the license, please [contact us](#).

4th December 2013

<http://hdl.handle.net/2440/81026>

Estimating a Markovian Epidemic Model Using Household Serial Interval Data from the Early Phase of an Epidemic

Andrew J. Black, Joshua V. Ross*

School of Mathematical Sciences, The University of Adelaide, Adelaide, Australia

Abstract

The *clinical serial interval* of an infectious disease is the time between date of symptom onset in an index case and the date of symptom onset in one of its secondary cases. It is a quantity which is commonly collected during a pandemic and is of fundamental importance to public health policy and mathematical modelling. In this paper we present a novel method for calculating the serial interval distribution for a Markovian model of household transmission dynamics. This allows the use of Bayesian MCMC methods, with explicit evaluation of the likelihood, to fit to serial interval data and infer parameters of the underlying model. We use simulated and real data to verify the accuracy of our methodology and illustrate the importance of accounting for household size. The output of our approach can be used to produce posterior distributions of population level epidemic characteristics.

Citation: Black AJ, Ross JV (2013) Estimating a Markovian Epidemic Model Using Household Serial Interval Data from the Early Phase of an Epidemic. PLoS ONE 8(8): e73420. doi:10.1371/journal.pone.0073420

Editor: Alessandro Vespignani, Northeastern University, United States of America

Received: March 6, 2013; **Accepted:** July 22, 2013; **Published:** August 30, 2013

Copyright: © 2013 Black, Ross. 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 author and source are credited.

Funding: This work was supported under the Australian Research Council's Discovery Projects funding scheme (project number DP110102893). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

* E-mail: joshua.ross@adelaide.edu.au

Introduction

A quantity which is commonly recorded during a pandemic is the *clinical serial interval*, defined as the time between date of symptom onset in an index case and the date of symptom onset in one of its secondary cases [1–3]. It was one of the main quantities recorded, at the level of households, during the 2009 H1N1 pandemic and subsequently used for understanding the dynamics of the pandemic [1,3–6]. Numerous studies have illuminated the critical dependence of disease dynamics and choice of control policy on this quantity through its relation to the generation time [7–10].

A common and simple way to analyse serial interval data is to fit it with a parametric distribution [4,11–13]. This approach allows an accurate calculation of the mean and possibly other moments. However, an obvious drawback of such an approach is that the estimate itself gives no understanding of the underlying mechanics, and hence it is difficult to make predictions with quantifiable confidence or to assess the impact of proposed control policies. This is because the serial interval is not a biological quantity in its own right but the convolution of the processes of transmission and incubation. This is further confounded by the fact that the time of infection is almost certainly unobservable, and because households are small, depletion of susceptibles has a large impact on the (stochastic) transmission process [14]. For these reasons, the only way to infer both epidemiological and dynamical quantities from serial interval data is by assuming and fitting a transmission model [15,16]. This approach not only provides an estimate of the serial interval distribution, but estimates a full mechanistic model which

may be used to make predictions and assess the impact of intervention policies [6].

A type of transmission model which has been growing in popularity, especially when considering household structure, are Markovian models [6,17–20]. In these it is assumed that there are two levels of mixing: strong mixing within a household and weaker mixing between households [17]. As the overall population is assumed to be large and randomly mixing, then during the early stages of an epidemic repeated introduction of infection into a single household is negligible. The assumption of only one introduction allows for deeper analysis of the model, and also allows for computationally-efficient methods to be developed for evaluating early-time quantities [17,21]; here it allows us to ignore the external infection rate, and use serial interval data to estimate the other parameters. Obviously during the mid-to-late stages of an outbreak, this assumption breaks down and hence more complex models are required, but for this study we confine ourselves to this common assumption. This early stage of an epidemic is important as we want to infer parameters which can then be used (with further assumptions) in population level models to assess possible interventions and inform public health policy.

In this paper we show how to fit a quite general Markovian household epidemic model using serial interval data. This is achieved by first explaining how the serial interval distribution can be calculated for this model, and hence used to derive exact likelihoods. We then use this for parameter inference via Bayesian Markov chain Monte Carlo (MCMC) methods. We investigate the accuracy of this methodology via simulation studies and illustrate its use with data previously studied from a household transmission study of seasonal influenza in Hong Kong [13]. Our investigations

Table 1. Within household dynamics.

Event	Transition	Rate
Infection	$(S, E_i) \rightarrow (S-1, E_{i+1})$	$\beta \frac{S \sum_{i=1}^k I_m}{(N-1)}$
Exposed progression, $(n = 1, \dots, j-1)$	$(E_n, E_{n+1}) \rightarrow (E_n-1, E_{n+1+1})$	$j\sigma E_n$
Start shedding	$(E_j, I_1) \rightarrow (E_j-1, I_{1+1})$	$j\sigma E_j$
Infection progression, $(m = 1, \dots, k-1)$	$(I_m, I_{m+1}) \rightarrow (I_m-1, I_{m+1+1})$	$k \gamma I_m$
Recovery	$I_k \rightarrow I_{k-1}$	$k \gamma I_k$

The transitions and associated rates which define the stochastic $SE(j)I(k)R$ model for the within-household dynamics.
doi:10.1371/journal.pone.0073420.t001

identify that household size has an appreciable impact on the serial interval distribution and that incorporating household size data into inference methods allows more accurate estimates of model parameters.

The advantages of our methodology are threefold: Firstly, it is fully stochastic and mechanistic – the former is vital given the average size of a household and the latter leads to greater understanding of the epidemic. Secondly, we can numerically solve the model, and hence calculate the serial interval distribution exactly to an arbitrary precision – there is no need for approximations by branching processes or for assumptions of independence in order to derive results. Thirdly, it is very computationally efficient. This means we can achieve the methodological ideal of full evaluation of the uncertainty in parameter estimates and derive accurate credible intervals for all results. Efficiency also allows for the potential inclusion of much more epidemiological detail in future models were more data available in such studies.

Methods

We assume a continuous time Markovian model for the dynamics of a disease within a household of size N . We use a general $SE(j)I(k)R$ model, where the exposed and infectious periods are split up into j and k phases so that each has an Erlang distribution with mean exposed and infectious periods $1/\sigma$ and $1/\gamma$, and variances $1/(j\sigma^2)$ and $1/(k\gamma^2)$, respectively [22,23].

Infection is assumed to be frequency dependent (but density-dependent transmission is no obstacle to the methodology we outline, and will be discussed later) with transmission parameter β [15,24]. The transition rates for this model are given in Table 1.

The model is specified by the matrix Q , which encodes the transition rates between different possible states of the household [6,18]. For the $SE(j)I(k)R$ model we consider, the total number of possible states is

$$\Psi = \frac{(1+j+k+N)!}{(1+j+k)! N!}, \tag{1}$$

hence this is also the dimension of Q . The element $Q_{\mu,v}$ is the rate of transition from state μ to v for $\mu \neq v$, where $\mu = 1, \dots, \Psi$ and $v = 1, \dots, \Psi$. $Q_{\mu,\mu} = -\sum_{v \neq \mu} Q_{\mu,v}$, is the negative of the total rate at which the system leaves state μ . The dynamics of the model are then given by the equation,

$$\frac{dp(t)}{dt} = p(t)Q, \tag{2}$$

where $p(t)$ is the probability vector with μ th entry the probability of the household being in state μ at time t [18]. As we are dealing with household models, the dimension of Q , given by Equation (1), is relatively small, so Equation (2) can be solved efficiently using matrix exponential methods [25]. Hence we can calculate $p(t)$ to an arbitrary precision, side-stepping the need for any type of potentially costly simulation.

To calculate the serial interval distribution we need to evaluate the probability of observing a secondary case in a given time

Table 2. Computational costs.

Household size, n	Time (s)	Effective Size
2	0.006	21
4	0.023	48
7	0.045	147

Average time taken to compute the likelihood for a household of size n . Other parameters $j=k=2$, $D_{max}=10$ and as given in Figure 1. A 2.5 GHz Intel core i5 machine running MATLAB was used for these timings. The Effective Size is the dimension of the Q matrix once the redundant states have been removed.
doi:10.1371/journal.pone.0073420.t002

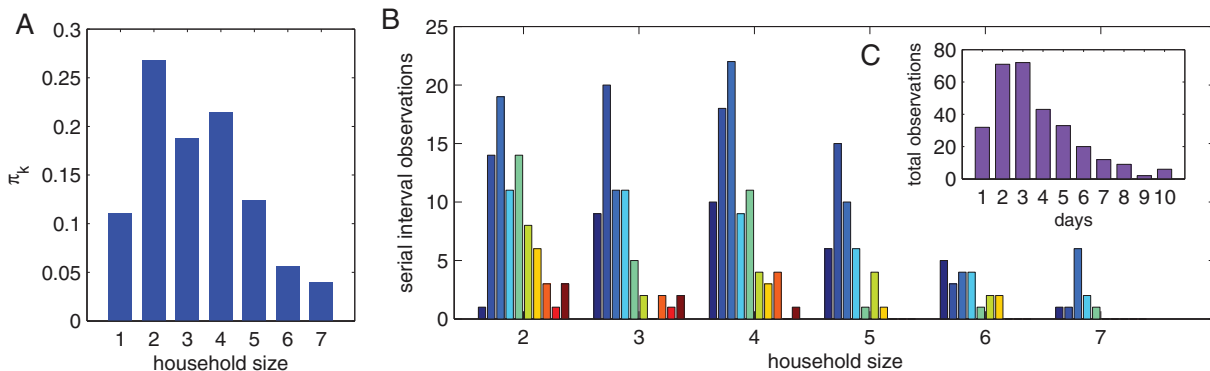


Figure 1. Generated serial interval distributions. (A) shows the size-biased distribution derived from USA 2011 census data. (B) shows 300 randomly generated serial interval observations, stratified by household size. (C) shows the same observations, but summed over all household sizes. These distributions are used in the next section to test the parameter inference methods. Parameters: $\beta=2$, $\sigma=1/4$, $\gamma=1/2$ and $j=k=2$.
doi:10.1371/journal.pone.0073420.g001

interval, given that we start with an index case at time $t=0$. To do this we first form the transition rate matrix Q corresponding to the $SE(j)I(k)R$ model for a given household size and set of parameters.

We assume that the appearance of symptoms coincides with entering the (first) infectious class [26]; in the later sections we discuss how Markovian models can be extended to weaken this assumption. The Q matrix is then modified so that states which correspond to a serial interval event – a second individual entering the first infectious class – are made absorbing. If we order the states of the system by $V=(E_1, E_2, \dots, E_j, I_1, I_2, \dots, I_k, R)$, then the set of absorbing states are

$$A = \left\{ \begin{array}{l} \{(E_1, E_2, \dots, E_j, 2, 0, 0, \dots, 0)\} \\ \{(E_1, E_2, \dots, E_j, 1, v_{k,1})\} \\ (0, 0, \dots, 0, 0, \dots, 1) \end{array} \right\} = \left\{ \begin{array}{l} B \\ (0, 0, \dots, 0, 0, \dots, 1) \end{array} \right\}, \quad (3)$$

where $v_{k,1}$ is the set of all vectors of length k , with a 1 in a single position and zeros elsewhere. The set of states B are those corresponding to serial interval events, while the last one is recovery with no further infection. This model explicitly takes into account that the second person to start showing symptoms might not have been the first to be infected, and hence evaluates the probabilities associated with the clinical serial interval.

For a household of size N , the initial state of the chain is set as $p(0)=(S=N-1, I_1=1)$. In doing this we are implicitly assuming that the first person to show symptoms is also the first person to introduce infection into the household. If we were considering asymptomatic individuals and/or multiple external infections then this might not be true. By numerically solving the dynamics we can then calculate the cumulative distribution function (cdf) of the serial interval, $F(t)$, by computing how much probability has flowed into the absorbing states by a given time. We then condition on the index case having created at least one secondary infection before recovering. The solution of the forward equation

giving the probability of being in a given state at time t is

$$p(t) = p(0) \exp(Qt). \quad (4)$$

Removing parts of the state space which are unreachable due to the new absorbing states can reduce the dimension of the matrix and speed up the evaluation of the matrix exponential. The cdf of the serial interval is then given by,

$$F(t) = \frac{1}{c} \sum_{s \in B} p_s(t), \quad (5)$$

where c is the probability that the index case infects at least one individual before recovering; note $c = 1 - \lim_{t \rightarrow \infty} \sum_{i \in \{A, B\}} p_i(t)$. The probability $1-c$ can be calculated simply by considering the sequences of events that would result in the individual going through k stages without infecting anybody. This then gives,

$$c = 1 - \left(\frac{k\gamma}{\beta + k\gamma} \right)^k \quad (6)$$

The serial interval probability mass function is formed by binning into days, as detailed in the next section.

Likelihood and MCMC algorithm

Given that we can compute the serial interval distribution for a given set of parameters to an arbitrary precision, calculating the likelihood for a given set of serial interval observations is relatively straightforward. Data on the serial interval is generally available at a daily resolution so we always work with a probability mass function binned into days. We used the following binning to calculate the probability of observing a secondary case on the i th day [4],

$$\rho_i = \frac{F(i+0.5) - F(i-0.5)}{F(D_{\max} + 0.5) - F(0.5)}, \quad i = 1, 2, \dots, D_{\max} \quad (7)$$

where $F(t)$ is the cdf and D_{\max} is the maximum range of observations. Given a set of index-secondary case observations, the likelihood of observing them is multinomial with probabilities ρ_i . If we have a number of household sizes then the likelihood is just the product of the likelihoods for each household size. MATLAB code to implement this procedure is provided via the Epistruct project [27].

The computational costs of calculating the likelihood are important. The dominant factor is the cost of evaluating the matrix exponential. The number of household sizes has the largest affect on the cost, and also larger households being relatively more expensive than smaller households due to their larger state spaces. Table 2 gives some average times to calculate the likelihood for individual household sizes using a 2.5 GHz Intel core i5 machine running MATLAB. The total average time to calculate the likelihood over $n=1, \dots, 7$ is 0.17 s using the same machine. The number of bins and the overall length of the distribution (D_{\max}) only have small effects on these timings as the EXPOKIT algorithm uses a variable step size [25]. The number of observations has no effect on the computational cost as these enter via a simple multinomial expression.

The method of inference was Bayesian MCMC [28]. To obtain samples from the posterior distribution we used a Metropolis-

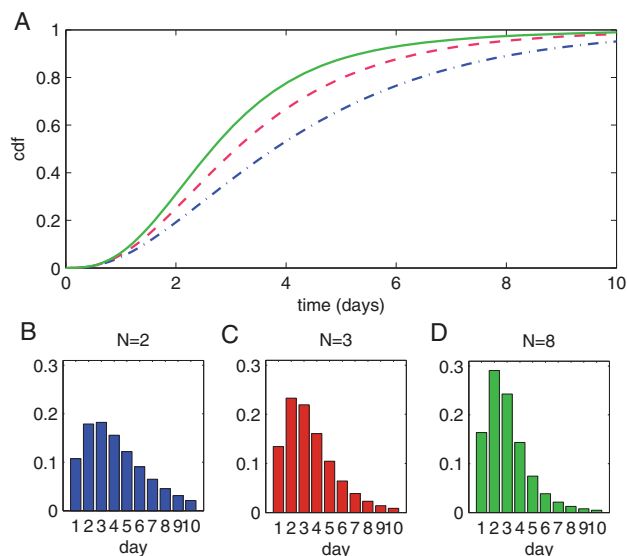


Figure 2. Theoretical serial interval distributions. Part A shows the serial interval cumulative distribution function for households of size $N=2, 3$ and 8 (dot-dashed, dashed and solid lines respectively). Parts B, C and D show the serial interval pmf (binned into days) derived from the corresponding cdfs. Parameter values: $\beta=2$, $\sigma=1/4$, $\gamma=1/2$ and $j=k=2$.

doi:10.1371/journal.pone.0073420.g002

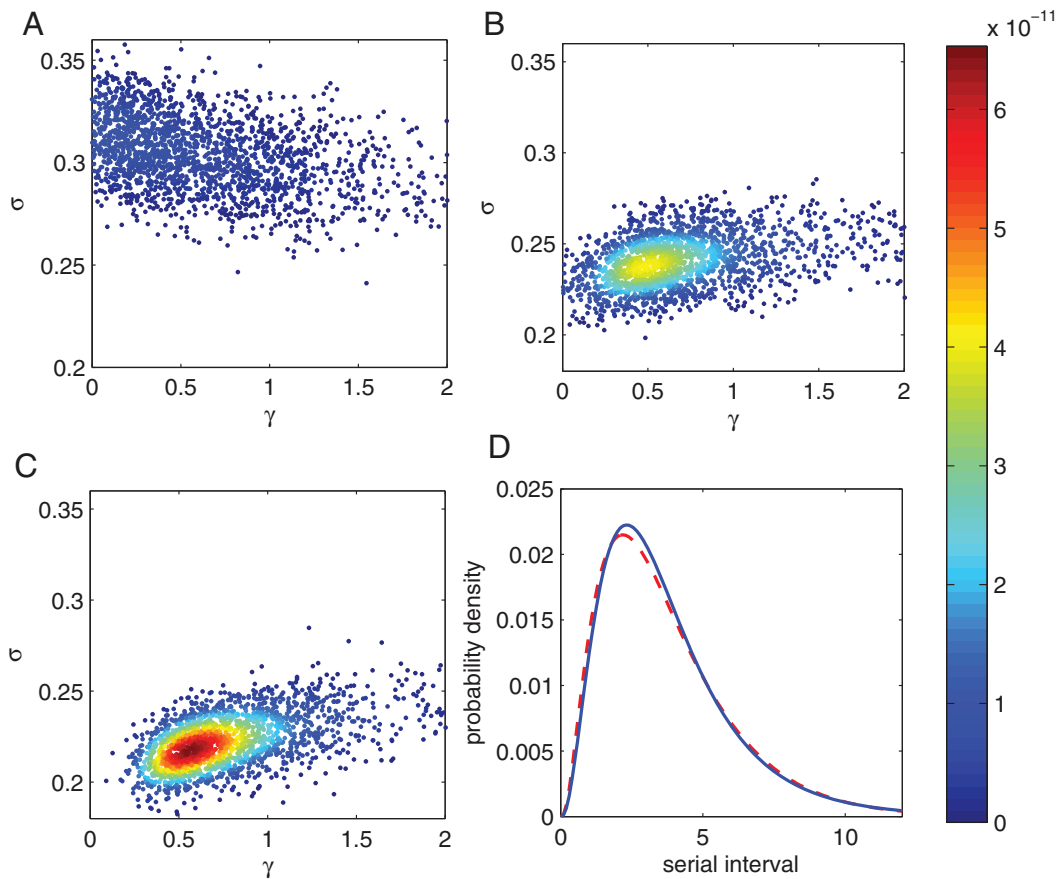


Figure 3. Parameter inference and predicted serial interval distributions. Plots A, B and C show 2×10^3 points from the posterior distributions for the parameters σ and γ assuming only a single household size, $N=2,3$ and 4 respectively. These are obtained from fitting to the distribution shown in Figure 1C. Points are coloured according to their likelihood with higher values assigned redder shades. All of these introduce a bias in the inferred parameters. Fixed parameters as in Figure 1. Part D shows the mean serial interval distribution for $N=2$ (dashed line) and $N=3$ (solid line). The distribution for $N=4$ is very close to that for $N=3$ so is not shown. True parameter values: $\beta=2$, $\sigma=1/4$ and $\gamma=1/2$. doi:10.1371/journal.pone.0073420.g003

Hastings algorithm with independent (truncated) Gaussian proposal densities. In all cases we assumed uniform priors on an interval from zero to an upper bound which depends upon the parameter. Burn-in time was 10^3 samples and the next 10^5 samples were taken, thinned by a factor of 10 to give 10^4 samples from the posterior; convergence was assessed via trace plots. The priors and trace plots for the individual runs are given in Appendix A of File S1.

Generating test data

To check the robustness of our method we generate a number of serial interval distributions with known household sizes and fixed parameters. We assume the early stages of an epidemic, so the distribution of infected household sizes will be approximately the size-biased distribution, $\{\pi_i\}$, where π_k is the probability of a randomly selected individual belonging to a household of size k [6,17]. This is given by

$$\pi_k = \frac{kh_k}{\sum_j jh_j}, \quad (8)$$

where h_k is the household size distribution for a given population. This provides a baseline case, obviously for household clinical trials a different distribution would be appropriate, but in any case

it would be explicitly known. Throughout this paper we use census data from U.S.A. 2011 for $h_k, k=1, \dots, 7$, which is shown in figure 1A.

The data is generated by first choosing a random number of household sizes (from 2 to 7) from the size-biased distribution. For each household, a serial interval observation is sampled according to the true distribution binned into days ($D_{\max}=10$). Figure 1B shows the simulated serial interval data stratified by household size. Figure 1C shows the simulated serial interval data summed over all household sizes.

Results

Effects of household size

Figure 2 shows how the serial interval distribution changes with household size, for sizes $N=2,3$ and 8 with frequency dependent mixing (β is held constant for different N). Larger households have higher probability of shorter serial intervals because there are more possibilities for who is the first individual to display symptoms. The change is greatest between $N=2$ and 3 , and decreases thereafter. This is because there is a trade-off between more people competing to show symptoms and the fact that these must have been infected later than the first person. As the household size increases the distribution therefore tends to a limiting case. As the variance of the exposed period decreases (j

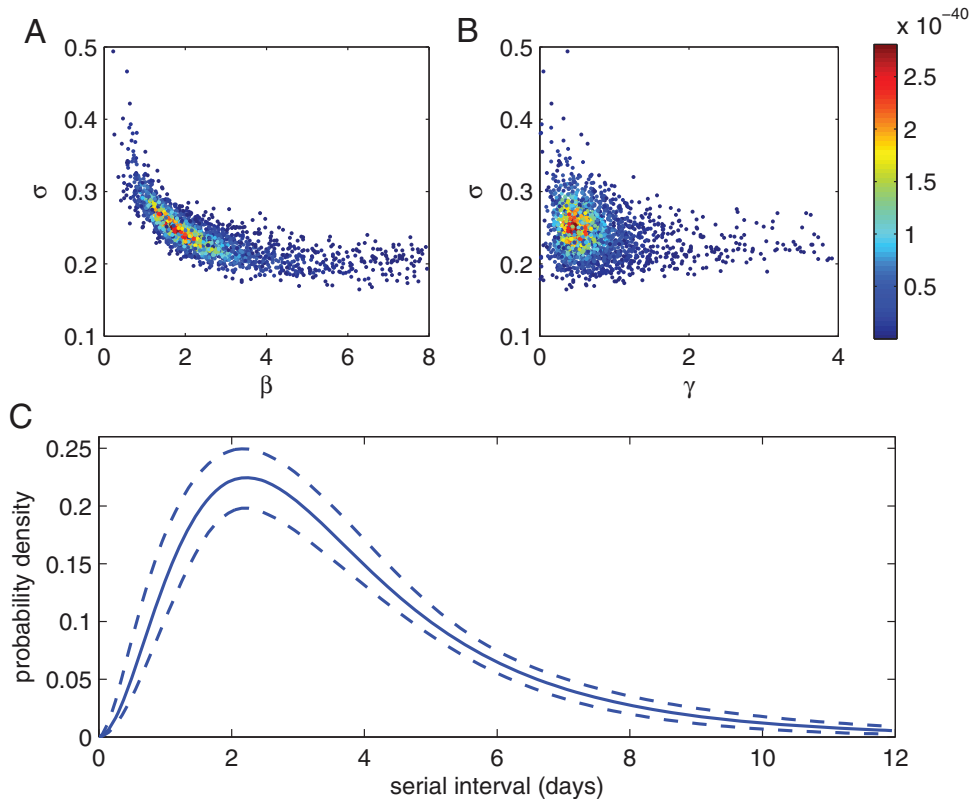


Figure 4. Inference of the serial interval distribution accounting for household size. (A) and (B) show 10^3 points from the posterior distribution projected along two different parameter axis. Points are coloured according to their likelihood with higher values assigned redder shades. (C) shows the mean serial interval distribution (solid line) and 95% credible intervals (dashed lines) obtained from 10^4 samples of the posterior, summed over all household sizes. True parameter values: $\beta=2$, $\sigma=1/4$ and $\gamma=1/2$. doi:10.1371/journal.pone.0073420.g004

increases) the serial interval also becomes more constant and the difference between the different sized households lessens. The variance of the infectious period (value of k) has only a small overall effect on the serial interval distribution, so henceforth we fix $k=2$ [6,29,30].

Inference with aggregated data

Here we report our findings when attempting to infer the posterior distribution for exposed period parameter σ and infectious period parameter γ by using the serial interval distribution assuming just a single household size—fixing all other parameters. In later sections we estimate all parameters, but here we are interested in quantifying the biases which can be introduced when using a single household size—effectively ignoring household size—to estimate the serial interval from data which has come from a population consisting of multiple household sizes. This situation often arises when trying to analyse aggregated data from the literature.

Figure 3 shows samples from the posterior distributions assuming three different (fixed) household sizes: $N=2,3$ and 4. The serial interval data used is that summed over all households, shown in Figure 1C, corresponding to a total of 300 serial interval observations. The case $N=2$ is biased by a large amount away from the true values, severely underestimating the infectious period parameter γ and overestimating the exposed period parameter σ . The $N=3$ case provides the best estimate of the parameters although there is still bias. Biases arising from using a model with $N>4$ grow larger, with σ underestimated and γ overestimated. The serial interval is most sensitive to the mean

exposed period, $1/\sigma$, and thus this is more accurately estimated. Although the parameter estimates from the three models are different, the estimated serial interval distributions corresponding to mean parameter estimates are all very similar (see Figure 3D), thus so are the mean serial intervals. The fit using $N=4$ is the best in terms of the mean likelihood.

Full inference from serial interval observations

We now look at estimation of all three variables: transmission parameter β , exposed period parameter σ and infectious period parameter γ , from the generated serial interval observations, given that we also know the household sizes for each observation, i.e. fitting to the data shown in Figure 1B. The variance of the exposed and infectious periods (parameters j and k) were held fixed. These can be inferred as well, but as noted earlier k cannot be inferred easily because the serial interval distribution is not very sensitive to it. In contrast the serial interval distribution is typically very sensitive to the variance in the exposed period (j) so in practice the actual value is almost always recovered. Figure 4 shows the posterior distributions along with the mean serial interval distribution with credible intervals. The MCMC algorithm for the full inference is appreciably slower than when using just a single household-size model, due to the higher dimension of the search space and need to calculate six individual serial interval distributions for each proposal.

To check the validity of our results we carried out sensitivity analysis. Specifically, we are interested in how the estimated posterior distribution depends on the number of observations available and how it can be skewed due to the random nature of

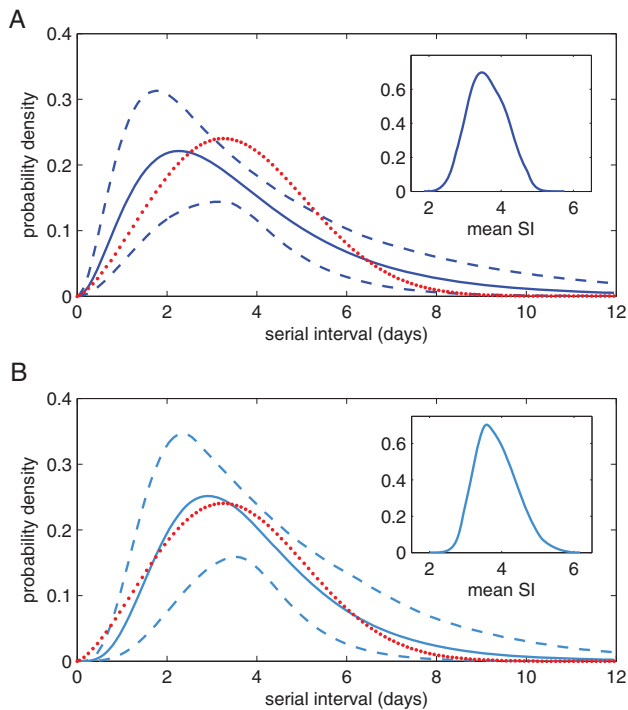


Figure 5. Estimated serial interval distributions. Solid lines depicts the mean and dashed lines the 95% credible intervals assuming $j=2$ (A) and $j=4$ (B). The dotted lines show the Weibull distribution estimated in the original analysis [13]. The insets show kernel density plots for the mean serial interval.
doi:10.1371/journal.pone.0073420.g005

the observations. To assess this we fit the full model to 8 sets of randomly generated serial interval distributions with 15, 50, 100, 200 and 300 data points respectively. The resulting posterior distributions are shown in Appendix B of File S1. The results of this show that exposed period parameter σ is found accurately most of the time, even for very small sample sizes. The other two parameters, transmission parameter β and infectious period parameter γ cannot be accurately determined until there are many more samples (typically at least 200). It is likely that we would need to include more of the later infection events within a household to resolve these parameters with more accuracy for smaller sample sizes.

It is also of interest to see how the estimates of parameters can be improved if one of the parameters is already established. We tested this by fixing the transmission parameter, β , and found the posterior distribution for the other two parameters (figures shown in Appendix C of File S1). This gives an improvement on the posterior for σ , but little improvement for γ . The serial interval distribution derived from this posterior has very similar credible intervals to that shown in Figure 4, so does not give an improved estimate for the mean serial interval.

Influenza in Hong Kong transmission study

We now use our model to estimate model parameters from a household study in Hong Kong [13]. In this study a Weibull model was fitted to clinical serial interval data corresponding to inter-pandemic influenza in Hong Kong during 2007. This was then used to estimate the mean serial interval with parametric bootstrapping to calculate confidence intervals [13]. Admittedly this has a very small sample size (only 14 observations from households of sizes $N=3$ to 5), but serves to illustrate the power of

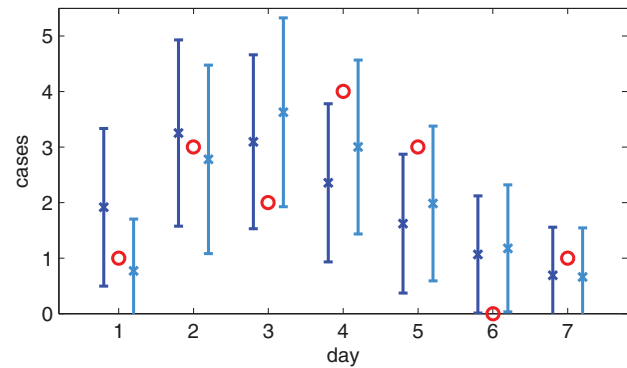


Figure 6. Expected number of serial interval observations of each duration compared with original data. Original data is summed over household size shown as circles (14 observations in total). The markers with bars show the mean and one standard deviations for the $j=2$ fit (x markers) and the $j=4$ fit (square markers).
doi:10.1371/journal.pone.0073420.g006

our method with real data. It is also the only study we have found which explicitly gives household size with serial interval observations. In the original study it was shown that external rates of infection had no impact on the serial interval estimate, so our model is appropriate to analyse this data.

To investigate the sensitivity to the variance of the exposed period we separately fitted two versions of the model with $j=2$ and 4. The higher value of j gives a more constant exposed period. As in the previous section we estimate the three parameters β , σ and γ and set $k=2$. Full details of the MCMC routine are given in Appendix A of File S1. Only the posterior for the exposed period parameter, σ , could be determined to within reasonable limits. Both values of j gave similar results: for $j=2$, $E(\sigma)=0.32$ (95% CI = 0.16–0.67), and for $j=4$, $E(\sigma)=0.32$ (95% CI = 0.18–0.68). The distributions for both the transmission parameter β and infectious period parameter γ were not well determined, but this is expected given the results of the sensitivity analysis in the previous section.

The estimated serial interval distributions and credible intervals are shown in Figure 5 for the two different values of j , along with kernel density plots for the mean serial intervals. In the original analysis a Weibull distribution was fitted [13] and is shown for comparison; the estimated mean serial interval was 3.6 days (95% confidence interval = 2.9–4.3). From the serial interval distributions we estimate the mean serial interval to be 3.6 days (95% CI = 2.6–4.6) assuming $j=2$ and 3.8 days (95% CI = 2.9–5.1) assuming $j=4$. The mean likelihood of the $j=4$ fit is approximately three times that of the $j=2$ fit. Figure 6 shows the expected number of serial interval observations of each duration and standard deviations for the two fits compared to the original data. For the $j=4$ case all the data lies within one standard deviation.

Discussion

The serial interval is relatively easy to observe and has been shown to be critically important for predicting disease dynamics and choosing control policies. For these reasons it is commonly recorded during the early stages of a pandemic. The difficulty arises when attempting to use the observations for modelling, or inference, because the serial interval is the convolution of two processes: infection and incubation, and the infection time is effectively unobservable.

In this paper we have provided methodology for parametrising a quite general Markovian model of household disease dynamics to serial interval data. Not only does this approach provide an estimate of the distribution of serial interval, but it also provides an estimate of a mechanistic model of the disease dynamics. This approach facilitates the prediction of disease dynamics and the assessment of alternative control options, of much importance in the early stages of disease invasion.

We have shown how the distribution of serial interval can be evaluated to arbitrary precision for our stochastic households model. Unlike stochastic simulation, which is computationally intensive and produces an estimate, our method is efficient and allows precise likelihood evaluation. Analytical results can be derived, using approximations in the cases $N > 2$ (see Appendix D of File S1), but in practice these offer no advantage over the numerical scheme due to the unwieldy nature of the expressions derived.

Our model allows us to quantify the effect of household size on the clinical serial interval (the time between first and second showing of symptoms, assuming that there is no asymptomatic infection and only a single introduction), hence identifying its importance for estimation. By fitting to serial interval distributions stratified by household size we can obtain accurate posterior distributions for all three of the basic parameters: transmission parameter β , exposed period parameter σ , and infectious period parameter γ . The parameter j , controlling the variance of the exposed period, can also be inferred, although we have not implemented this within the MCMC scheme. The serial interval distribution is relatively insensitive to the parameter k , controlling the variance in the infectious period, so we have not attempted to infer this and have held it constant. If full time series of symptomatic events are available then our method is potentially wasteful because we do not use the later events. Our methods can be extended to inference of full time series and it is likely that this is required to get better estimates on the parameters β and γ . Such a project is currently under way.

We have shown the effectiveness of this scheme for estimating parameters from simulated data as well as data from a Hong Kong influenza study [13]. Despite the small sample size we could still infer meaningful estimates for the exposed period and serial interval distribution, consistent with the earlier study. This demonstrates that the methodology reliably produces estimates as would be obtained via traditional parametric fitting, but has the added benefit of producing estimates of parameters for our stochastic, mechanistic model of disease dynamics. Of course, one must be careful in using household quantities to make population level predictions [9]; to do this we typically need to make more assumptions about population level mixing and transmission. In related work on antiviral effectiveness [6] we have used this method with a simpler model to effectively estimate the exposed period parameter σ and infectious period parameter γ from a larger influenza serial interval dataset [4]. Although this dataset was larger, the data was not stratified by household size, so we had to use a mean household size in our estimation. This then allowed us to evaluate posterior distributions for population level quantities such as the household basic reproductive number, R_* , and early growth rate [17,21,31].

References

1. First Few Hundred (FF100) Project (28 May 2009). Health Protection Agency, Health Protection Scotland, Communicable Disease Surveillance Centre Northern Ireland, and National Public Health Service for Wales; United Kingdom. Available: http://www.hpa.org.uk/webc/HPAwebFile/HPAweb_C/1257260453727. Accessed 2013 Aug 5.

The serial interval is also important because of its relation to the generation time which can be used to relate the Malthusian early growth rate, r , and the basic reproductive ratio, R_0 [8,9,14,23]. Usually it is assumed that these two distributions have the same mean, but in general their distributions will be different [32]. The actual generation time distribution can be derived for our model in a similar way to the serial interval distribution. Briefly, one would change the initial condition of the Markov chain to $E_1 = 1$ and make a different set of states absorbing. Once the joint posterior distribution for the parameters is inferred from the serial interval data, we can use it to compute the generation time distribution.

Whilst our model is quite general, there exists a number of features which may be required for particular diseases, populations and data sets which would require modification of our approach. For example, we have not explicitly accounted for external infection and co-primary cases, varying infectiousness with stage of infection, or symptoms that do not coincide with the commencement of infectiousness. It possible to extend the model to account for these features, and the method we have outlined will also need to be modified slightly to accommodate these extensions. We note that in all cases extra parameters will require estimation. We are currently developing and testing such frameworks. However, the model we have explicitly analysed herein is of much interest in infectious disease modelling, and the method we have detailed will facilitate its use in the early stages of disease invasion, of much interest to public health policy.

Here we have shown that household size has a significant impact on the serial interval, and that including this data improves estimates. Throughout we have assumed frequency-dependent transmission, as appears to be most appropriate for influenza in households [15], but one would expect the differences to be exacerbated by density-dependent transmission – not only do larger households have more individuals competing to display symptoms first, but the transmission rate would also be larger for the same household configuration which further reduces the serial interval. Household size is typically recorded alongside the serial interval, so our method simply proposes a way to make appropriate use of this routinely collected data; an approach which has the benefit of producing posterior distributions of parameters corresponding to a fully-mechanistic model of the disease dynamics.

Supporting Information

File S1 This file contains Statistical (MCMC) details; Sensitivity analysis of the full model; Full inference while holding β fixed; and, Some analytic results.
(PDF)

Acknowledgments

We thank the referees for their helpful comments.

Author Contributions

Conceived and designed the experiments: AB JR. Performed the experiments: AB JR. Analyzed the data: AB JR. Contributed reagents/materials/analysis tools: AB JR. Wrote the paper: AB JR.

4. Donnelly CA, Finelli L, Cauchemez S, Olsen SJ, Doshi S, et al. (2011) Serial intervals and the temporal distribution of secondary infections within households of 2009 pandemic influenza A(H1N1): implications for influenza control recommendations. *Clin Infect Dis* 52(S): S123–S130.
5. Lau LLH, Nishiura H, Kelly H, Ip DKM, Leung GM, et al. (2012) Household transmission of 2009 pandemic influenza A(H1N1): a systematic review and meta-analysis. *Epidemiology* 23: 531–542.
6. Black AJ, House T, Keeling MJ, Ross JV (2013) Epidemiological consequences of household-based antiviral prophylaxis for pandemic influenza. *J R Soc Interface* 10: 20121019.
7. Germann TC, Kadau K, Longini IM, Macken CA (2006) Mitigation strategies for pandemic influenza in the United States. *PNAS* 103: 5935–5940.
8. Ferguson NM, Cummings DAT, Cauchemez S, Fraser C, Riley S, et al. (2005) Strategies for containing an emerging influenza pandemic in Southeast Asia. *Nature* 437: 209–214.
9. Wallinga J, Lipsitch M (2007) How generation intervals shape the relationship between growth rates and reproductive numbers. *Proc R Soc B* 274: 599–604.
10. Griffen JT, Garske T, Ghani AC (2011) Joint estimation of the basic reproduction number and generation time parameters for infectious disease outbreaks. *Biostatistics* 12: 303–312.
11. Lipsitch M, Cohen T, Cooper B, Robins JM, Ma S, et al. (2003) Transmission dynamics and control of severe acute respiratory syndrome. *Science* 300: 1966–1970.
12. Lessler J, Reich NG, Cummings DAT (2009) Outbreak of 2009 pandemic influenza A(H1N1) at a New York city school. *N Engl J Med* 361: 2628–2636.
13. Cowling BJ, Fang VJ, Riley S, Peiris JSM, Leung GM (2009) Estimation of the serial interval of influenza. *Epidemiology* 20: 344–347.
14. Scalia Tomba G, Svensson A, Asikainen T, Giesecke J (2010) Some model based considerations on observing generation times for communicable diseases. *Math Biosci* 223: 24–31.
15. Cauchemez S, Carrat F, Viboud C, Valleron AJ, Boëlle PY (2004) A Bayesian MCMC approach to study transmission of influenza: application to household longitudinal data. *Stat Med* 23: 3469–3487.
16. Cauchemez S, Donnelly CA, Reed C, Ghani AC, Fraser C, et al. (2009) Household transmission of 2009 influenza A (H1N1) virus in the united states. *N Engl J Med* 361: 2619–2627.
17. Ball F, Mollison D, Scalia-Tomba G (1997) Epidemics with two levels of mixing. *Ann App Prob* 1: 46–89.
18. Keeling MJ, Ross JV (2008) On methods for studying stochastic disease dynamics. *J R Soc Interface* 5: 171–181.
19. Lipsitch M, Abdullahi O, D'Amour A (2012) Estimating rates of carriage acquisition and clearance and competitive ability for pneumococcal serotypes in Kenya with a Markov transition model. *Epidemiology* 23: 510–519.
20. Pitzer VE, Basta NE (2012) Linking data and models: The importance of statistical analyses to inform models for the transmission dynamics of infections. *Epidemiology* 23: 520–522.
21. Ross JV, House T, Keeling MJ (2010) Calculation of disease dynamics in a population of households. *PLoS ONE* 3: e9666.
22. Lloyd A (2001) Destabilization of epidemic models with the inclusion of realistic distributions. *Proc R Soc B* 268: 985–993.
23. Wearing HJ, Rohani P, Keeling MJ (2005) Appropriate models for the management of infectious diseases. *PLoS Med* 2: e174.
24. Keeling M, Rohani P (2008) Modeling infectious diseases in humans and animals. Princeton University Press.
25. Sidje RB (1998) Expokit: A software package for computing matrix exponentials. *ACM Trans Math Softw* 24: 130–156.
26. Carrat F, Vergu E, Ferguson NM (2008) Time lines of infection and disease in human influenza: A review of volunteer challenge studies. *Am J Epidemiol* 167: 775–785.
27. Epistruct (2013). MATLAB routines for epidemic modelling and inference in structured populations. Available: <http://sourceforge.net/projects/epistruct/>. Accessed 2013 Aug 5.
28. Gilks WR, Richardson S, Spiegelhalter DJ (1995) Markov Chain Monte Carlo in Practice. Chapman and Hall/CRC.
29. House T, Baguelin M, van Hoek AJ, Flasche S, White P, et al. (2011) Modelling the impact of local reactive school closures on critical care provision during an influenza pandemic. *Proc R Soc B* 278: 2753–2760.
30. Baguelin M, van Hoek AJ, Jit M, Flasche S, White PJ, et al. (2010) Vaccination against pandemic influenza A/H1N1v in England: A real-time economic evaluation. *Vaccine* 28: 2370–2384.
31. Ross JV (2011) Invasion of infectious diseases in finite, homogeneous populations. *J Theor Biol* 289: 83–89.
32. Svensson Å (2007) A note on generation times in epidemic models. *Math Biosci* 208: 300–311.