

Evaluating Particle Filtering Methods for Complex Multi-Pathogen Disease Systems

Xi Meng¹ and Nicholas G. Reich¹

¹Department of Biostatistics, University of Massachusetts, Amherst

Abstract As a popular framework of likelihood-based inference, particle filtering has found applications in a wide range of settings. We evaluated this method in a complex disease system where immunological interactions between multiple pathogens give rise to challenges in drawing inference for parameters. In particular, we adapt a published epidemiological model with 80 states and 13 parameters that represents a biologically accurate model of immunological interactions between four serotypes of dengue fever. Specifically, we are interested in estimating two parameters of particular biological interest, cross-protective duration and antibody-dependent enhancement. A standard application of particle filtering and multiple iterated filtering for simulated data from this system leads to results that do not converge, even when many parameters are assumed to be fixed at the true value. A central challenge is that the likelihood surface cannot be estimated precisely. Therefore, we developed and implemented a particle filtering approach that adaptively focuses on regions of the parameter space with high estimated likelihood values. This approach enabled us to characterize the profile likelihood surfaces for these parameters. Coupled with a method to construct confidence intervals in the presence of noisy likelihood estimates, we constructed maximum likelihood estimator and confidence intervals. We tested the approach on a set of simulated data as well as real dengue case-count data. In datasets simulated from realistic parameter values, we were able to estimate accurately the impact of enhancement (avg. bias = 0.16; 95% CI coverage probability = 0.9). However, our ability to estimate the cross-protection parameter was limited (avg. bias = 0.4 year; 95% CI coverage probability = 0.4). Additionally, the method showed little sensitivity to misspecification of the parameters assumed to be known. Our analyses quantify the challenges of implementing a full-scale analysis using these methods for complex disease transmission models.

1 Introduction

Dengue is a mosquito-transmitted virus that is the cause for dengue fever (DF) and dengue hemorrhagic fever (DHF). It has been estimated that dengue virus infects 50 million people worldwide per year. Approximately 2.5 billion people live in countries with dengue epidemic; most of them are in tropic areas [1]. Dengue virus has four closely related but distinct serotypes. Each of the four serotypes is capable of causing acute infectious diseases of DF and DHF.

Various interacting factors among serotypes, such as cross protection and susceptibility enhancement, make transmission of multi-serotype infectious disease a complicated process. Infections with any of the four serotypes of dengue virus can provide long-term protective immunity to the same serotype of virus, but only provide short-term protective immunity to a different serotype of virus. Secondary infections with a serotype different from a previous infection are at high risk of resulting in more serious diseases, due to a consequence called antibody-dependent enhancement (ADE) [2,3].

Research has shown evidence that disease models fit data better if they take these serotype interactions into account [4-6]. Reich et. al. (2013) analyzed serotype-specific dengue infections data in Thailand and found quantitative evidence that dengue infections provide short-term immune interactions among serotypes. They compared models including and excluding immunological interactions, and found that models including cross protection perform better in predicting dengue incidence. Previously the study of Adams et.al. (2006)

showed models with cross-protective immunity match better with the observed epidemic pattern. Wearing et.al. (2006) also found it was necessary to include cross-immunity as well as antibody-dependent enhancement to simulate dengue incidence with similar features to observed data.

As the importance of taking serotype interactions into account is becoming increasingly evident, there is growing interest in drawing inference about interaction parameters among different serotypes. However, it is challenging to establish a mathematical disease progression model because immunological interactions between multiple serotypes make disease transmission a complicated process, and infectious disease data only provide partial information about the disease dynamics.

A class of likelihood-based inference methods proposed by King and Ionides, particle filtering [7-11] and iterated filtering [12-14], has found applications in a range of disciplines including infectious disease epidemiology. Broadly speaking, these methods are used to approximate posterior distributions of hidden states in state-space models with sets of weighted particles. These methods have advantages over traditional techniques in drawing inference from infectious disease settings that are strongly nonlinear and contain latent variables of the underlying transmission process [7,9,15]. First, for state-space models with two stochastic processes of state process and observations, these methods provide great tools to track the states of the dynamic system. Additionally, they outperform standard alternatives such as Markov chain Monte Carlo methods in high-dimensional systems [10]. Moreover, they are computationally-efficient by finding an approximate model representation without solving an exact solution of the full likelihood, more efficient than Kalman filters [15].

Particle filtering and its extension technique to maximize likelihood, iterated filtering, have been successfully applied in studies of cholera dynamics [12,17], pertussis immunity [18], measles [19], malaria [20], and poliovirus transmission [21], and Ebola [22]. Shrestha et.al. (2011 & 2013) have shown correct inference of interactions in two-pathogen systems could be obtained by an inference framework associated with particle filtering [23,24]. However no studies have showed the feasibility of this likelihood-based method in a 4-serotype dengue system, and its feasibility when misspecification exists in some epidemiological information.

In this paper, we assessed the capacity of the likelihood-based method of particle filtering to draw inference for unknown parameters associated with serotype interactions. We applied this likelihood-based inference framework to an established multi-serotype epidemiological model for a 4-serotype dengue system, with the goal of drawing inference about two immunological interaction parameters, duration of cross-protective immunity, and the level of antibody dependent enhancement. We used an adaptive particle filtering approach across changing two-dimensional grids of parameter space focusing on higher likelihood estimates. Coupled with a profile likelihood method We constructed adjusted confidence intervals and MLE from likelihoods with high uncertainty. We evaluated the performance of this method in a small simulation study, and assessed the sensitivity of this statistical inference method when there is presence of misspecification in initial states cases and transmission rate. Additionally, we applied this method on a real data set of serotype-specific dengue incidence in Bangkok, Thailand. This approach enables us to find the confidence interval for parameters of interest, and shows some promise for estimating the effects of immunological interactions between different pathogens, especially for the level of ADE. While this work shows some successes, it also quantified the challenges of implementing a full-scale analysis using particle filtering approach for a complex disease transmission model.

The remainder of the paper is organized as follows: we describe inference method based on particle filtering and disease transmission modeling in Section 2; Section 3 contains simulation studies; an application to the Bangkok dengue incidence is given in Section 4, and we conclude with discussions in Section 5.

2 Methods

2.1 Overview of particle filtering

The idea of particle filtering [7-11] is to use a number of draws of “particles” to represent the distribution of model states. The set of particles propagate forward the disease transmission model in time according to a set of ordinary differential equations. And then they assimilate the next observation record to refine the estimate. Observations are used to give weight to each particle based on its likelihoods. Particles closer to observations are given higher weight while those not close to observations are given lower weight. The number of particles used in each run is associated with the estimation accuracy. Increasing the number of

particles will lead to lower variability in estimates but result in higher computation burden [25]. The process of particle filtering will output log likelihood estimates for each possible parameter set, but will not provide maximum likelihood estimation. The algorithm of particle filtering we use is *Algorithm 1* in King (2015) [25].

Maximum likelihood estimation via iterated filtering (MIF) [12-14] is a method to obtain an optimized model with maximum likelihood estimator. MIF is based on particle filtering. The difference is that parameters in particle filtering are time-invariant, while model parameters in MIF are time varying with a random walk that decrease across each iteration. After carrying out multiple particle filtering operations, when the random walk intensity approaches zero, the likelihood will converge to the maximum. Unlike particle filtering, MIF provides maximum likelihood estimation instead of calculating likelihoods for fixed parameter combinations. The algorithm of MIF is according to Algorithm 3 in King (2015) [25].

We use the framework of partially observed Markov process (POMP) [25] to construct a transmission model and conduct inference tests. POMP is a convenient tool for dealing with time series data that provides a flexible framework for constructing nonlinear, non-Gaussian state-space models that are connected by a Bayesian network. Particle filtering and MIF along with some other modern statistical methods are implemented in the R package *POMP*.

A general POMP model has two essential components: the process $\{X_t\}$ and the observations $\{Y_t\}$ made at times t_1, \dots, t_n . Observations are assumed to be conditionally independent given the process. The process model is generated according to some equations that in our case represent a biologically accurate disease transmission process. Our observation model, which is the observed monthly cases of serotype i at time t , $y_{i,t}$, is assumed to follow a Poisson distribution,

$$y_{i,t} \sim \text{Poisson}(\lambda_{i,t})$$

where ρ is reporting rate linking the process and observation model, $\lambda_{i,t} = \rho x_{i,t}$.

The conditional transition probability density function is written as $f_{X_t|X_{t-1}}(x_t|x_{t-1}; \Theta)$, and the measurement probability density function is represented by $f_{Y_t|X_t}(y_t|x_t; \Theta)$, where Θ is the parameter vector, $\{Y_t\} = \{y_1, \dots, y_n\}$ represent the sequence of observation data at times t_1, \dots, t_n , and $\{X_t\} = \{x_0, \dots, x_n\}$ represent the sequence of process data at times t_0, \dots, t_n . The log likelihood function given parameters is sum of conditional log likelihoods,

$$\ell(\Theta) = \sum_{t=1}^n \ell_{t|t-1}(\Theta)$$

where $\ell_{t|t-1}(\Theta) = \log f_{Y_t|Y_{1:t-1}}(y_t|y_{1:t-1}; \Theta)$. In our POMP model,

$$\ell_{t|t-1}(\Theta) = \log \int f_{Y_t|X_t}(y_t|x_t; \Theta) f_{X_t|Y_{1:t-1}}(x_t|y_{1:t-1}; \Theta) dx_t$$

Since $\ell(\Theta)$ has no closed form, we turn to particle filtering and use representations of $f_{X_t|Y_{1:t-1}}(x_t|y_{1:t-1}; \Theta)$ to obtain an estimated $\hat{\ell}(\Theta)$ to approximate $\ell(\Theta)$.

The likelihoods obtained from particle filtering method contain large Monte Carlo error. Because of high computation amount required, it is unrealistic to lower the Monte Carlo error to a few log likelihood units in this complex model. In order to correct for the likelihood estimates noise, we used the profile likelihood method proposed by Ionides, et al (2016) to construct an adjusted confidence interval with Monte Carlo uncertainty taken into account.

2.2 Transmission model

2.2.1 Model outline

To model the characters of a four-serotype dengue transmission system, we apply a Susceptible-Infected-Convalescent-Recovered (SICR) model generated from a two-serotype infectious model by Shrestha (2011). The whole host population is split into four compartments of infectious status: susceptible, infected, convalescent, and recovered, denoted by S, I, C, R , respectively. Individuals are assumed to enter the system susceptible to all four serotypes. The model equations are shown in Section 2.2.2. In this model we assume subjects have equal susceptibility to all four dengue pathogens prior to getting infections. And we assume all dengue virus serotypes are equal in terms of cross-protection duration, susceptibility enhancement, and transmission rate.

simulation parameter	assumed value	estimated or fixed
cross-protective duration $1/\delta$	1 day to 2 years	estimated
antibody dependent enhancement χ	0.5 to 1.5	estimated
population size N	10,000,000	fixed
mortality rate μ	0.02/year	fixed
birth rate μ	0.02/year	fixed
transmission rate β	70/year	fixed
infectious period $1/\gamma$	2 weeks	fixed
interaction during infection ϕ	0	fixed
interaction during convalescent ζ	0	fixed
environmental reservoir ω	5e-7	fixed
environmental stochastic SD in transmission rate τ_{sd}	0.005	fixed
reporting rate of primary infection ρ_1	1	fixed
reporting rate of secondary infection ρ_2	1	fixed

Table 1: Model Parameters

As shown in Figure 1, the transmission process for each serotype follows a Susceptible-Infected-Convalescent-Recovered development. Individuals follow the paths downwards and go through four infections with all four serotypes. We assume there are serotype interactions for recovered individuals, but no cross-interactions for currently infected or convalescent individuals; in other words, the interaction parameter during infection or convalescent equal to zero.

The model has 80 initial states and 13 parameters (Table 1), 2 of which, cross-protective immune duration and antibody-dependent susceptibility enhancement, are to be estimated using the likelihood-based inference method.

2.2.2 Deterministic equations

Let X denotes the number of people at a specific infectious status (termed a compartment). X_0 is the number of people susceptible to all four serotypes and immune to none. The subscripts of I, C, R represent infected, convalescent, and recovered, respectively. The subscripts of i, j, k, l represent the four serotypes of dengue virus. For example, X_{I_j, R_i} represents those infected by serotype j and recovered from serotype i . $X_{I_i, R.}$ represents those infected by serotype i and recovered from other strains. Refer to Table 1 for other parameter notations such as infection rate, transmission rate, and so on.

The deterministic equations are as following:

- (1) Initial Status (susceptible to 4 serotypes)

$$\frac{dX_0}{dt} = \mu N - (\lambda_i + \lambda_j + \lambda_k + \lambda_l + \mu)X_0$$

- (2) 1st infection (infected by a serotype i)

$$\begin{aligned} \frac{dX_{Ii}}{dt} &= \lambda_i X_0 - (\gamma + \mu)X_{Ii} \\ \frac{dX_{Ci}}{dt} &= \gamma X_{Ii} - (\delta + \mu)X_{Ci} \\ \frac{dX_{Ri}}{dt} &= \delta X_{Ci} - [\chi(\lambda_j + \lambda_k + \lambda_l) + \mu]X_{Ri} \end{aligned}$$

- (3) 2nd infection (recovered from serotype i , infected by another serotype j)

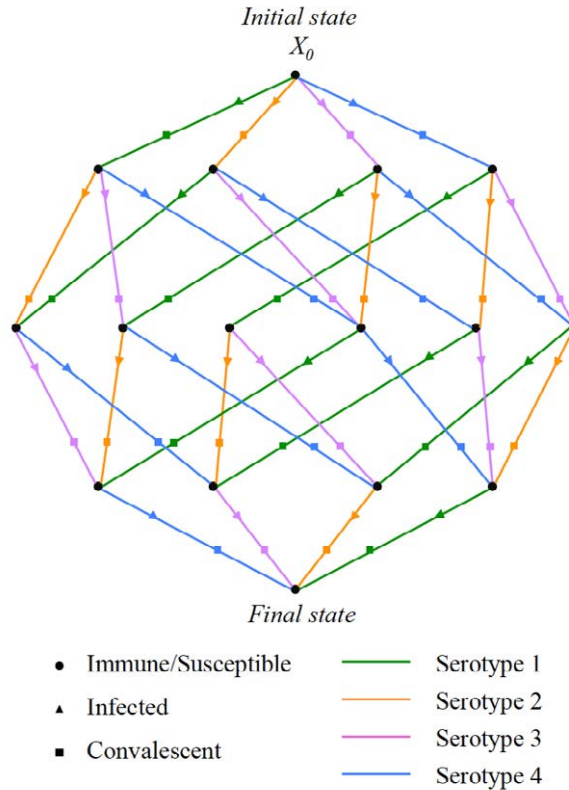


Figure 1: Model Schematics. The 80 infectious state are represented by nodes of different shapes. The node at the bottom represents the state of being immune to all four serotypes due to infection with all of them. The transmission process starts from the top node, where susceptible to all four serotypes and immune to none, and follows the paths downwards until the bottom node, where immune to all four serotypes and susceptible.

$$\begin{aligned} \frac{dX_{Ij,Ri}}{dt} &= \chi\lambda_j X_{Ri} - (\gamma + \mu)X_{Ij,Ri} \\ \frac{dX_{Cj,Ri}}{dt} &= \gamma X_{Ij,Ri} - (\delta + \mu)X_{Cj,Ri} \\ \frac{dX_{Rij}}{dt} &= \delta(X_{Cj,Ri} + X_{Ci,Rj}) - [\chi(\lambda_k + \lambda_l) + \mu]X_{Rij} \end{aligned}$$

(4) 3rd infection (recovered from serotypes i, j , infected by a third serotype k)

$$\begin{aligned} \frac{dX_{Ik,Rij}}{dt} &= \chi\lambda_k X_{Rij} - (\gamma + \mu)X_{Ik,Rij} \\ \frac{dX_{Ck,Rij}}{dt} &= \gamma X_{Ik,Rij} - (\delta + \mu)X_{Ck,Rij} \\ \frac{dX_{Rijk}}{dt} &= \delta(X_{Ck,Rij} + X_{Ci,Rjk}) - (\chi\lambda_l + \mu)X_{Rijk} \end{aligned}$$

(5) 4th infection (recovered from serotypes i, j, k , infected by a fourth serotype l)

$$\begin{aligned} \frac{dX_{Il,Rijk}}{dt} &= \chi\lambda_l X_{Rijk} - (\gamma + \mu)X_{Il,Rijk} \\ \frac{dX_{Cl,Rijk}}{dt} &= \gamma X_{Il,Rijk} - (\delta + \mu)X_{Cl,Rijk} \\ \frac{dX_{Rijkl}}{dt} &= \delta(X_{Cl,Rijk} + X_{Ci,Rjkl} + X_{Cj,Rikl} + X_{Ck,Rijl}) - \mu X_{Rijkl} \end{aligned}$$

2.2.3 Stochasticity

We introduced stochasticity by multiplying the transmission rate by a noise term. The infection rates are defined as

$$\lambda_s = \left[\frac{\beta}{N} (X_{Is} + X_{Is,R.}) + \omega \right] \frac{dw}{dt}, \quad \forall s \in \{i, j, k, l\}$$

where dw/dt represents white noise with gamma distribution, $dw \sim \text{Gamma}(\frac{dt}{\tau_{sd}^2}, \tau_{sd}^2)$.

We assume the transitions between model states to be a random process, specifically, an independent multinomial process with expected rates determined by above equations.

For example, the number of infected individuals coming out of the completely susceptible population is

$$y \sim \text{Multinomial}(X_0, \pi)$$

where $\pi = (\frac{\lambda_i}{p}, \frac{\lambda_j}{p}, \frac{\lambda_k}{p}, \frac{\lambda_l}{p}, \frac{\mu}{p})$, $p = \lambda_i + \lambda_j + \lambda_k + \lambda_l + \mu$.

2.3 Particle filtering inference for the model

To estimate cross-protective immune duration and susceptibility enhancement, we form a 2-dimensional parameter grid and use particle filtering to obtain likelihoods given parameters over possible combinations in the parameter space. In this complicated model, we found a standard application of maximum likelihood estimation via iterated filtering lead to results that do not converge, even when many parameters were assumed to be fixed at the true value. As a result, we used an adaptive particle filtering approach across increasingly focused two-dimensional grids of parameter space in order to estimate two parameters of interest.

Coupled with the method of profile likelihood correcting Monte Carlo error [26], we constructed maximum likelihood estimators and confidence intervals for each simulated dataset. The likelihood profile for one parameter is constructed by fitting a lowess smoother with weights depending on proximity to the approximated quadratic peak [26]. For each value of one parameter, we use the highest three average likelihood over the other parameter, like Figure 2e and 2f. The maximum likelihood estimator is taken as the maximum of the quadratic approximation.

The particle filtering algorithm proceeds in the following steps:

Step 1: Form an initial 8-by-8 2-dimensional grid of immune period($1/\delta$) and susceptibility enhancement(χ) covering a large range of reasonable values from 1 day to 4 years for δ , and from 0.1 to 5 for χ . It is evenly spaced in both parameters.

Step 2: Apply particle filtering to obtain likelihoods given parameters over varying possible combinations in the parameter space. There are 64 combinations of parameters in total.

Step 3: Based on the likelihoods obtained from particle filtering, we select a narrower parameter space including the grid points that have a log-likelihood within the highest 5%. The new vector contains 8 evenly spaced elements, ranges from $\theta_1 - I/2$ to $\theta_2 + I/2$, where θ_1 and θ_2 are the two bounds, I represents the interval width of the last grid. Ideally we will have a narrower parameter space with less interval and covering the true value.

Step 4: Conduct further particle filtering with higher computation intensity on the narrower parameter space. We used 20,000 particles in each iteration, 3 filters in the first two iterations and 5 filters in the later iterations.

Step 5: Repeat steps 2-4 for several times and find the parameter vector with the maximum likelihood and the 95% adjusted confidence interval for the two focal parameters using the methods of profile likelihood.

2.4 Simulation study design

By way of a proof-of-concept, we assume all parameters to be known and generate serotype-specific monthly infections, shown in Figure 7a. Our simulated data consists of number of infections with four different serotypes of a multi-serotype infectious disease, observed at each month in 40 years. On the basis of simulated infectious data, we attempt to infer back the model parameters. We integrate the time series data and a transmission model within a likelihood-based inference framework implemented in R-package *pomp*. We assume that initial conditions for all infectious states and other epidemiological parameters are known and only attempt to infer the parameters pertaining to the interactions among serotypes, susceptibility enhancement, and the duration of cross-protection.

We simulate 90 data sets of serotype-specific monthly infectious data for a length of 40 years from the epidemiological model described previously by varying only the parameters to be estimated. We generate 9 scenarios of data sets using different combinations of varying cross-protection duration ($1/\delta$) and the level of ADE (χ). For each scenario, we generate 10 datasets of the same parameter set. Given the computational time needed to analyze a single dataset (150 CPU hours), analyzing more simulated datasets is not feasible.

We tried to obtain a maximum likelihood estimator with multiple iterated filtering by starting from the parameter value with highest log-likelihood from the last particle filtering run and using 20 replicates and 50 iterations. However, very few of the 90 data sets converge to a maximum likelihood estimator. In most cases the estimate variation increases with number of iterations with no sign of convergence. Thus we turn to use the particle filtering iteration method applied on gradually decreasing parameter grid in the direction of higher likelihood.

3 Methods evaluation: simulation study

3.1 General performance on simulated data

The ability to draw inference depends on parameters. In datasets simulated from realistic parameter values, we observed an average bias of 0.02 year and 0.2 in estimating the duration of cross-protection between serotypes and the level of antibody dependent enhancement, respectively. Parameter estimates for each data set in Figure 3 are obtained by the method of profile confidence intervals in the presence of noisy likelihoods. Figure 3a shows the estimations for short-term protective duration grouped by level of ADE. Figure 3b shows the estimations for the level of ADE grouped by short-term protective duration. Variabilities in estimating both parameters are considerable, while the estimation of antibody dependent enhancement χ showed less variability than that of immune duration $1/\delta$.

3.2 Performance with longer time length

Further simulation study showed improvement in inference ability with more data, especially for estimating the level of antibody dependent enhancement χ .

In the previous simulation study, we assessed the method on dataset of length of 40 years. In order to improve the estimation accuracy and get a better understanding of the method performance on this complex model with serotype interactions, one approach could be generating longer time series which contain more information about the transmission dynamics. We propose three different conduction procedures with different amount of computation intensity and test them on two groups of simulated data, 40-year length and 100-year length. In particular, we simulated 10 data sets for each time length, with the same realistic parameters, temporary protection of 2 years and ADE level of 1.5. Each dataset is estimated by three procedures.

Ideally the length of confidence interval and bias (Figure 4) will decrease with more iterations and maintain a high rate of capturing the true value (Figure 5) at the same time. According to Figure 4 and Figure 5, we obtained more accurate estimates and higher coverage probabilities in simulated data with longer time length.

3.3 Performance with higher computational intensity

Results also showed increasing the amount of computation intensity by using large number of particles and filters doesn't have large impact on inference ability.

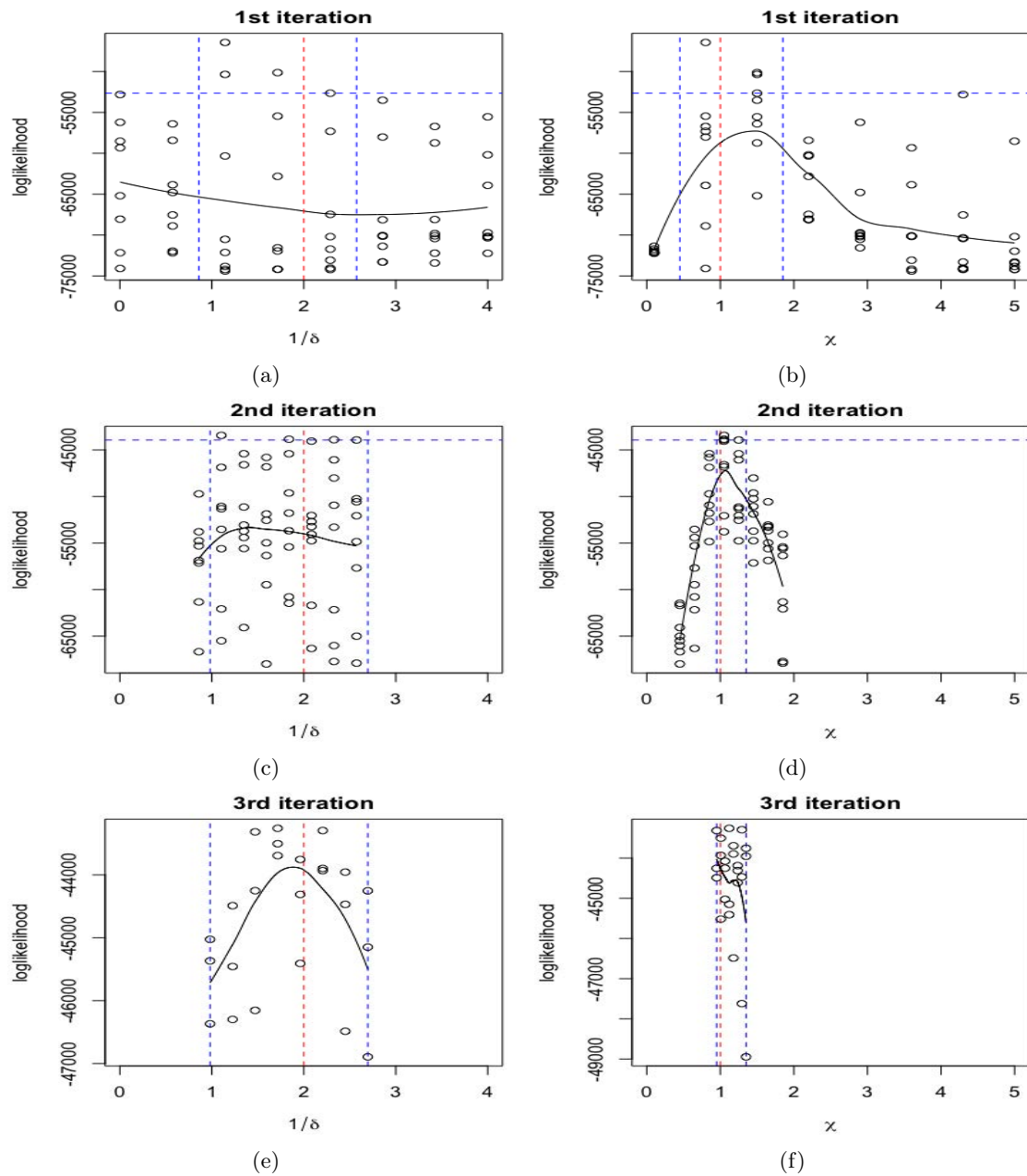
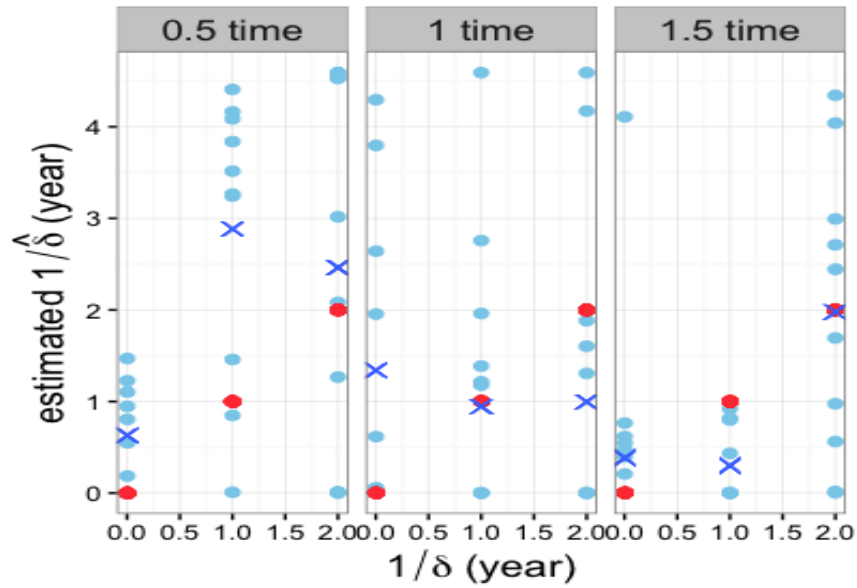
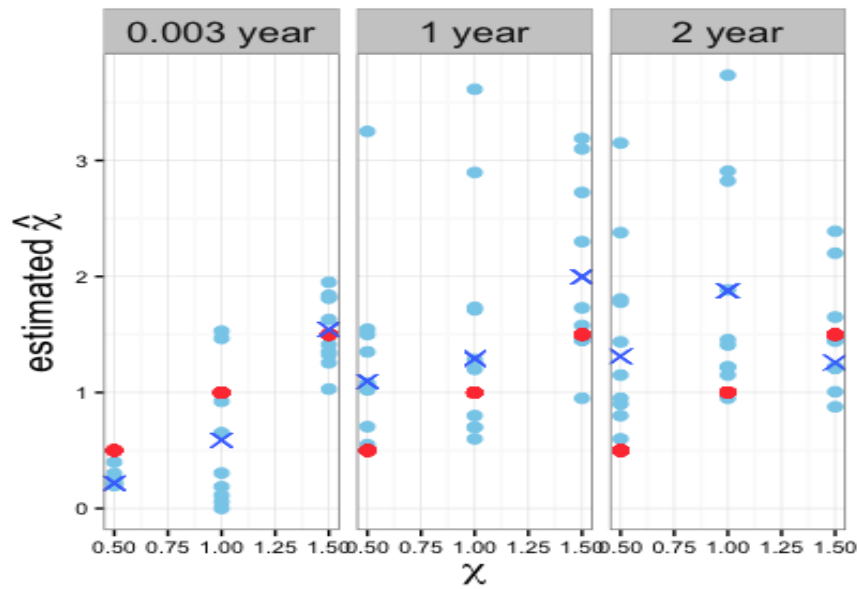


Figure 2: Likelihood profile. Left panel are likelihood profiles for parameter $1/\delta$ from each of the three particle filtering iterations, right panel are those for parameter χ . True parameter values are represented by red lines. The focused parameter space are generated by the parameter values corresponding to the highest 5% log likelihoods. The 95% confidence intervals in the third iteration are generated by the profile likelihood method [26] based on results from the last iteration.



(a) Faceted plot of $1/\delta$ vs $\hat{1/\delta}$ by χ



(b) Faceted plot of χ vs $\hat{\chi}$ by $1/\delta$

Figure 3: Estimation results from 9 scenarios of simulated data sets. Blue dots represent maximum likelihood estimators for each data set, red dots represent real values used to simulate the data. Blue X's are average estimations for each scenario of data sets.

Previously we assessed the method using 20,000 particles in each particle filtering (PF) iteration. Besides using longer time series, another approach that might improve the estimation accuracy could be increasing computational intensity. Thus, we propose three different conduction procedures each using different numbers of particles and filters, and test them on two groups of simulated data.

For the 1st procedure, we use 20,000 particles in each iteration, 3 replicating filters in the first two iterations, and 5 filters in the last iteration. For the 2nd procedure, we use 20,000 particles, 3 filters in the first iteration; use 40,000 particles, 3 filters in the second iteration; and use 40,000 particles, 5 filters in the last iteration. For the 3rd procedure, we use 20,000 particles, 3 filters in the first iteration; use 40,000 particles, 3 filters in the second iteration; and use 80,000 particles, 5 filters in the last iteration.

By comparing length of confidence interval and bias (Figure 4), and 95% coverage probability (Figure 5) at each iteration step, we find computation intensity has little impact on estimating accuracy. Although procedures with larger computation intensity give more accurate estimates in cross-protective duration, it cannot compensate the computation cost.

In a word, this inference method allows us to estimate accurately the impact of susceptibility enhancement. Both length of intervals and bias decrease with iterations and the estimate bias goes down to essentially zero after the third iteration. And coverage probabilities maintain high rates, especially in data with 100-year length. However the ability to estimate cross-protective duration is limited. In the same datasets with 100-year length that simulated by realistic parameter values, the average bias for level of susceptibility enhancement is 0.16, 95% CI coverage probability = 0.9; the average bias for level of cross-protective duration is 0.4 year, 95% CI coverage probability = 0.4.

3.4 Sensitivity analyses

The likelihood-based methods are moderately robust to misspecification of other parameters. Estimation based on parameters with uncertainty shows no difference from that without uncertainty.

Previously, we assumed that epidemiology data rather than interactions were an unbiased reflection of the true value. In reality some of these parameters are obtained through reporting surveys, and thus can be biased and need to be estimated. To examine the influence of uncertainty magnitude on the estimation accuracy using this likelihood-based framework, we conducted sensitivity analyses. Analyses were carried out by modifying our assumptions about known parameters and testing the capacity of this approach in the presence of uncertainty. We conducted the analysis in two situations, one is misspecification of one parameter, transmission rate; the other is misspecification of several parameters, initial conditions.

Examination is performed on the 10 datasets simulated from realistic parameter values of $1/\delta = 2yr$ and $\chi = 1.5$.

(1) Misspecification in one parameter. We added different degrees of misspecification to transmission rate. When constructing transmission models we used transmission rate values that are some degree (from 10% higher to 20% lower) away from its true value. Then we evaluated the performance in drawing inference for the two interaction parameters.

(2) Misspecification in a set of parameters. We evaluated the sensitivity of this likelihood-based inference method in the situation where misspecification exists in initial conditions of infectious states. To generate misspecification in initial states, we assumed each initial condition X_0^i as a Normal distribution with expected mean of the true value, and added a random standard deviation sd^i so that 95% of the distribution is within 10%-30% of the mean.

$$X_0^i \sim Normal(X_0^i, sd^i)$$

where $sd^i = cvX_0^i$, and cv is the coefficient of variance.

Both of the sensitivity analyses are shown in Figure 6, with a comparison to the group without uncertainty in other parameters. Estimation with uncertainty shows no difference from that without uncertainty. There is a considerable amount of variation in protective duration, which is consistent with previous results. Relative to the average estimate without uncertainty, misspecification introduces a small amount of bias.

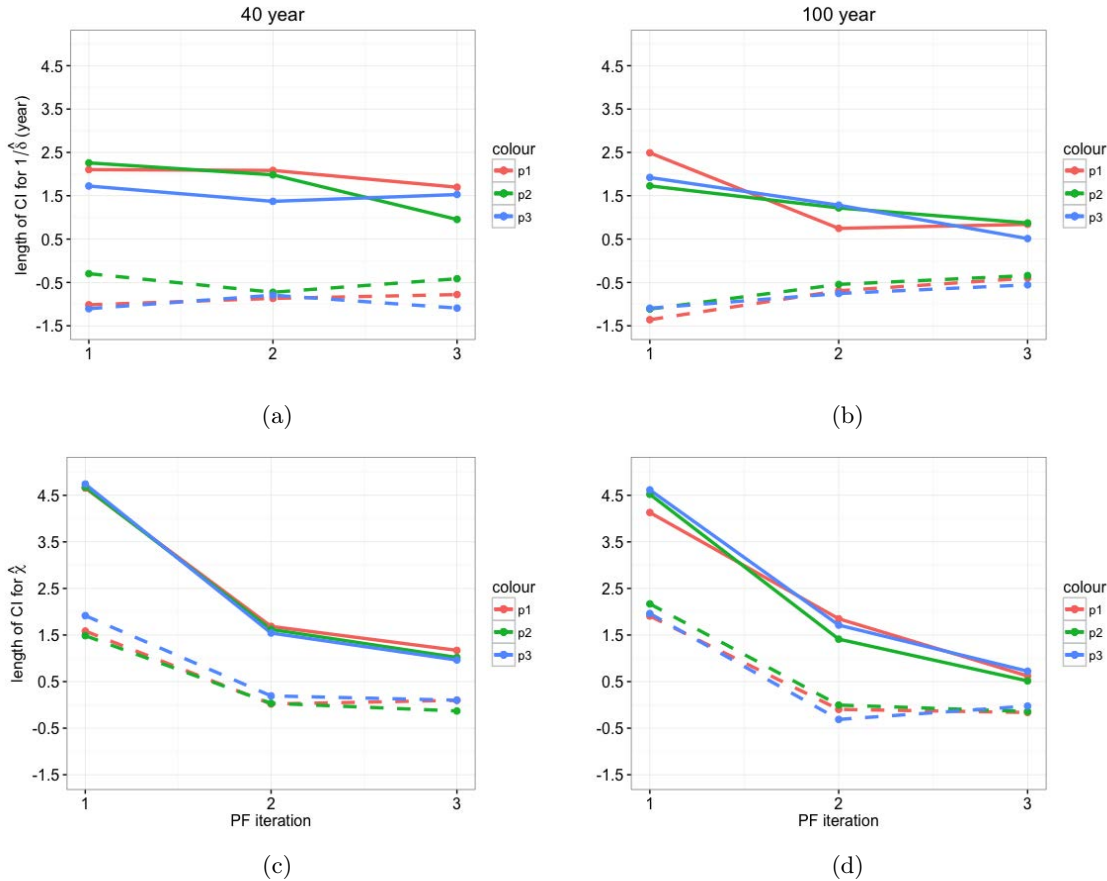


Figure 4: Average length of confidence interval and average bias by different procedures at each iteration step. Each dataset is estimated using three procedures, represented by different colors. Dashed lines are estimated bias for different procedures. 4a and 4b shows the confidence interval length and the estimate bias for temporary protection in the 40-year data and that in the 100-year data. 4c and 4d shows the interval length and bias for ADE in the 40-year data and that in the 100-year data.

4 Application: Interaction Estimation in Bangkok Dengue Cases

We apply our method to the data of monthly case-count in Bangkok from 1973 to 1999. The dataset consists of serotype-specific laboratory confirmed dengue cases of 27 years in Bangkok, Thailand. The data was reported by the Queen Sirikit National Institute of Children’s Health, which is a children’s healthcare facility serving as a reference hospital for dengue in Bangkok.

Since the data has a low reporting rate, there are many zero observations that make the estimation difficult to perform. To solve it, we reconstruct true estimated case-count by adding a small value and dividing by the reporting rate at the corresponding time points. The reproduced time series of serotype-specific monthly case counts for 27 years assuming reporting rate equals 100% is shown in Figure 7b.

Before applying the inference approach on this time series, we need to estimate the 80 initial conditions. We proceed as follows:

- (1) Reproduce the case data with reporting rate of 100% from reporting case c_t and reporting rate ρ_t , $(c_t + 0.1) \times \rho_t$.
- (2) For each grid point in the parameter space, simulate the model and get 10 sets of state values at 10 different time point.
- (3) For each initial guess, quantify the fit of the model to the real data using the synthetic log-likelihood based upon several summary statistics. Calculate the average synthetic log-likelihood for each grid point.

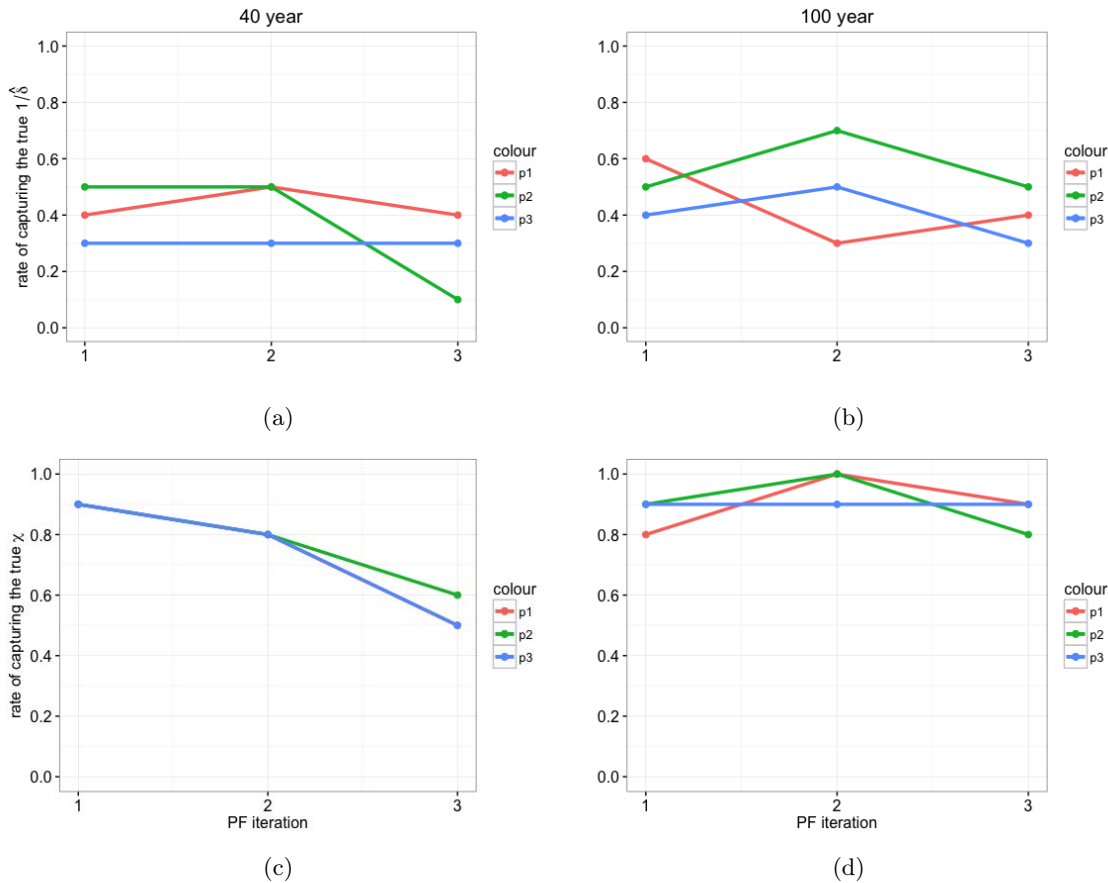


Figure 5: Coverage probabilities by different procedure at each iteration step. Each dataset is estimated using three procedures, represented by different colors. 5a and 5b shows coverage rate of true values of temporary protection in the 40-year data and that in the 100-year data. 5c and 5d shows rate of capturing true values of ADE in the 40-year data and that in the 100-year data. In 5c procedure 1 and 2 overlap together.

- (4) Select the initial conditions estimated with the highest synthetic likelihood.
- (5) Proceed with the particle filtering approach stated above.

The likelihood profiles for the parameters of interest after each iteration is shown in Figure 8. Methods for choosing focused parameter space and constructing confidence interval are same with those used in simulations. The figure shows a wide 95% CI range for protective duration covering almost the whole range of reasonable values, which is consistent with the low identifiability for protective duration that we found in the simulation studies. Confidence interval for susceptible enhancement is from 0.3 to 2.3.

5 Discussions

Interactions in multi-pathogen infectious disease transmission give rise to significant challenges to researchers seeking to understand the nature of epidemiologic processes. The impact of immunological interactions among dengue serotypes introduces difficulty in constructing an accurate disease model and drawing inference about the parameters. The likelihood-based particle filtering framework is among the few approaches that could be applied in these situations.

Possessing unique advantages over other techniques, this likelihood-based inference method is an increasingly popular approach to solve estimation problems in epidemiology. Particle filtering is a convenient tool for solving dynamic problems in nonlinear state-space models with two stochastic components of process and

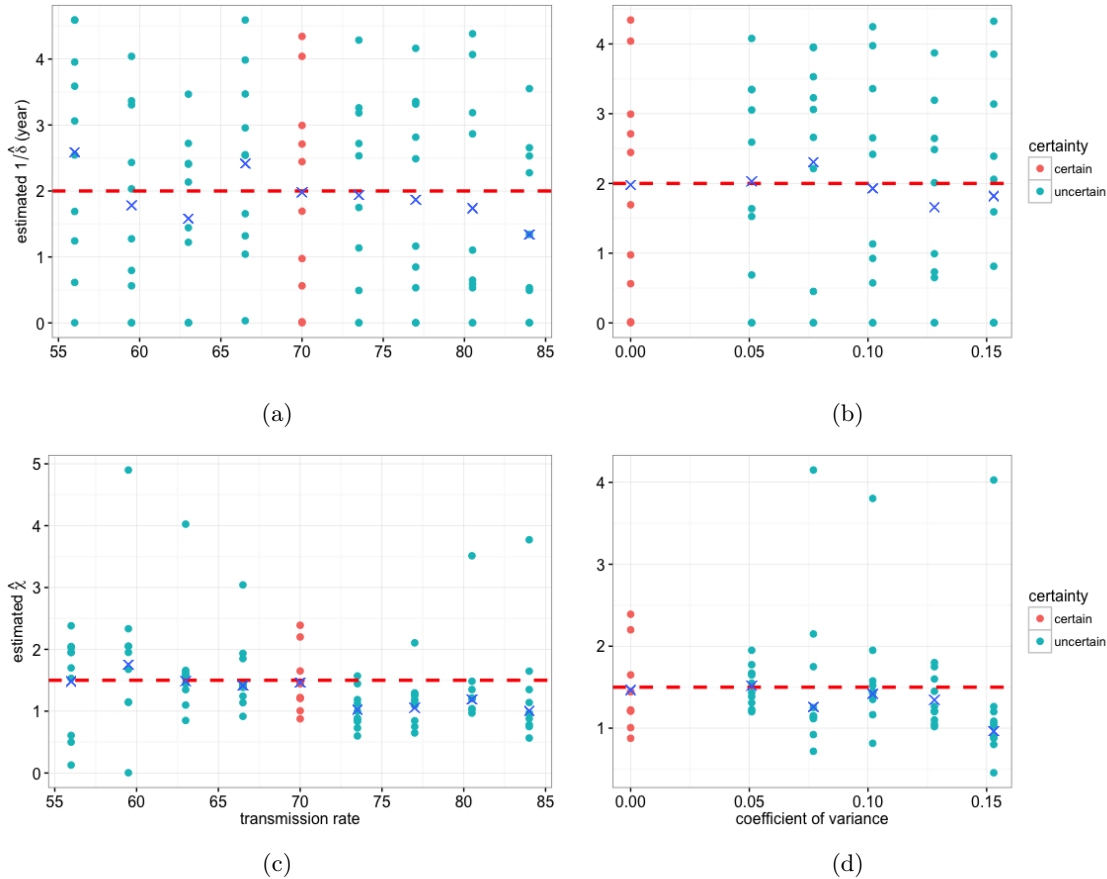
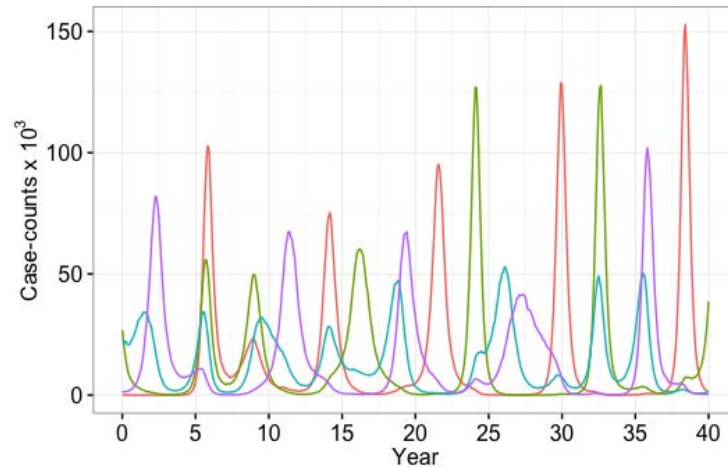


Figure 6: Sensitivity Analysis. 6a and 6c represent estimated temporary protection and ADE with misspecification in transmission rate, 6b and 6d are estimated temporary protection and ADE with misspecification in initial conditions. Blue dots represent maximum likelihood estimators with uncertainty in other parameters, red dots represent estimates without uncertainty, blue X's are average estimates, red dotted line indicates the true value.

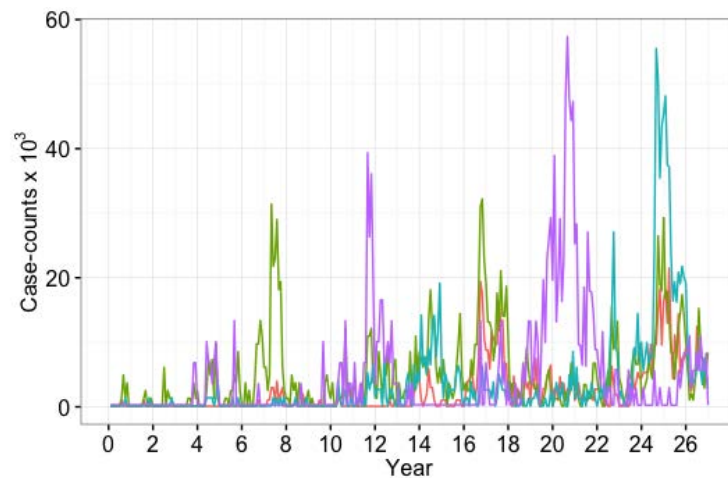
observations. Traditional methods like MCMC can also be applied in this kind of problems, but it is not as convenient as particle filtering in dealing with high-dimensional problems [10]. Another advantage of the particle filtering is that it does not need to evaluate the transition density of the latent Markov process. It seeks to find an approximate system representation without solving an exact solution of the model. In this way it outperforms a similar technique called Kalman filtering [9,15]. Furthermore, it is easy to implement and suitable for parallel implementation [16]

Apart from these theoretical advantages, in a complex model as the 4-serotype-4-infection model with 80 states and 13 parameters in this paper, the inference ability for particle filtering is not extremely good. Likelihood shows some promise as a basis for estimating the effects of immunological serotype interactions. Results suggest that level of ADE is more capable to be accurately estimated by particle filtering method, while short-term protection is not well identifiable. On the identifying ability, MIF provides the same information as particle filtering. We have tried different parameter setting for running MIF, including increasing particles or changing cooling fractions; but it did not converge to a unique identifiable maximum likelihood estimator. One reason might be the incidence data provides more direct information on level of susceptibility enhancement than short-term cross-protective immune duration.

By modifying our assumptions about other known parameters such as transmission rate and initial states cases, we conducted sensitivity analyses to test the robustness of the results in the presence of uncertainty. Results show little impact of misspecification in drawing inference for immunological interactions. The likelihood-based inference performs similarly in situations with misspecification and without misspecification



(a) Simulated serotype-specific dengue cases



(b) Serotype-specific dengue cases in Bangkok

Figure 7: Simulated data and real data

in known parameters. By evaluating the sensitivity of uncertainty in initial conditions, we found evidence for a simpler inference process by reducing the burden of estimating a large number of initial conditions. The ability of insensitive to misspecification in initial conditions has significance in reality. Initial conditions could be considered known parameters instead of unknowns when estimating interactions.

One limitation for this inference technique is the computational intensity. The whole process using the least computation expensive procedure described previously for one single data set takes about 150 CPU hours. And due to the extensive calculation time for the whole process, we only ran a small simulation study of 10 datasets for each of 9 scenarios, and we implemented particle filtering on a sparse parameter grid over a large area. Results might be more convincing if we conduct analysis in a larger simulation study. And confidence interval construction based on asymptotic likelihood ratio test would be more appropriate if implementing on a dense parameter grid. Another limitation is the high standard variation in likelihoods among different filters for the same parameter set, which makes the estimation less identifiable and also introduces challenges for constructing reliable confidence intervals. Currently we used the method of profile likelihood quadratic approximation to correct the uncertainty in log likelihoods, but there remains some problems. Our study has two unknown parameters and each parameter set yields several likelihood evaluations from different filters. To construct a likelihood profile one needs to find the maximum likelihood over the other nuisance parameters. The current method of constructing confidence intervals accounts for the error from data and the

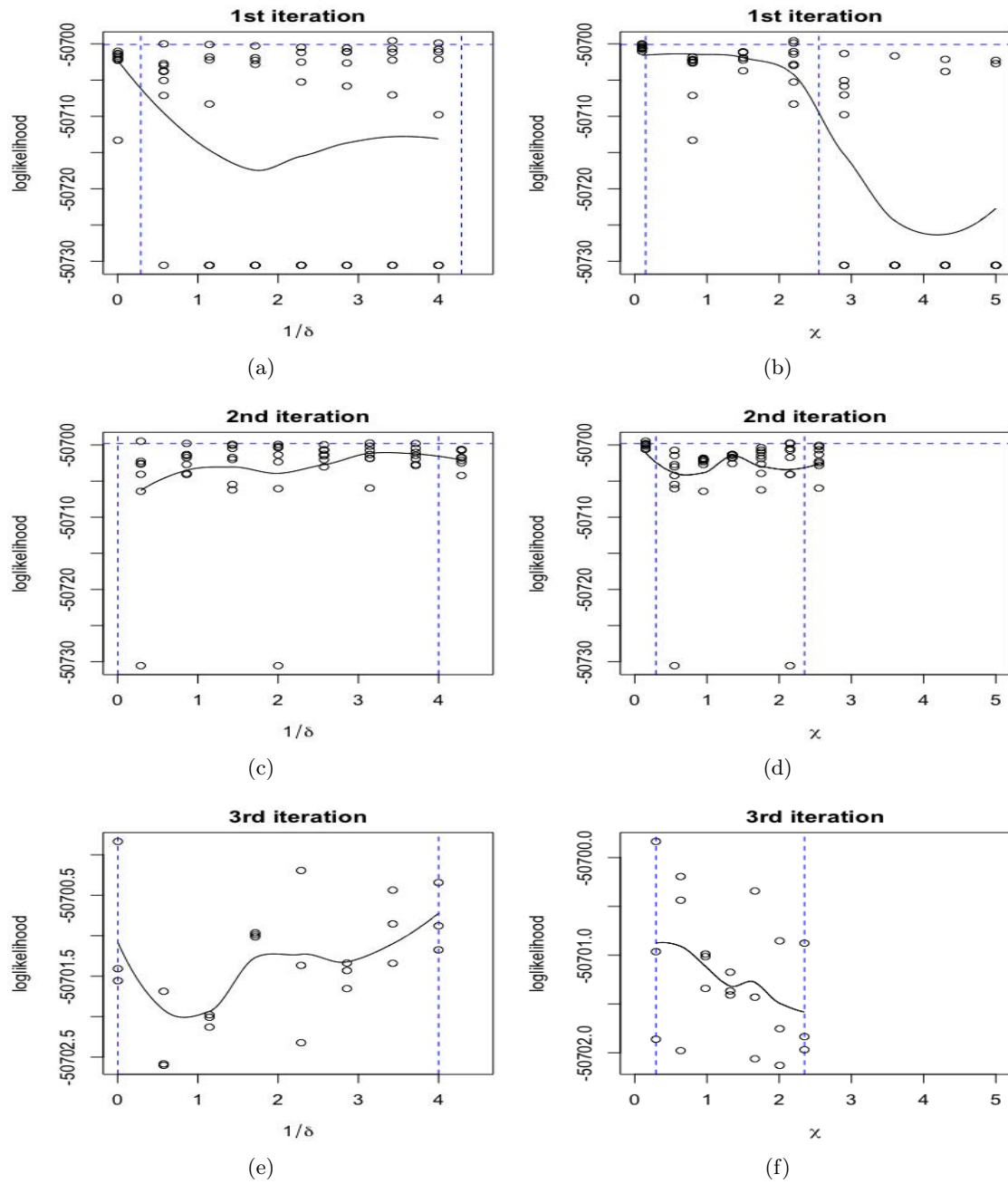


Figure 8: Likelihood profile of real data. Left panel are likelihood profiles for parameter $1/\delta$ from each of the three particle filtering iterations, right panel are those for parameter χ . Blue lines in 8e and 8f correspond to estimated confidence intervals for the two parameters.

error from Monte Carlo estimator, but doesn't account for the error from estimating the maximum likelihood of the other nuisance parameters.

While an inferential approach in this complex system based on particle filtering has strong theoretical backing and has proved useful in simpler systems, our analyses highlighted and quantified some of the challenges of implementing a full-scale analysis using particle filtering for a four-serotype dengue transmission model. Our study is one of the first feasibility analyses for this likelihood-based inference method in complex models. Although a simple model does not help to fully understand the nature of immunological serotype interactions, our model is too detailed that it introduces uncertainty and brings down the power of inference ability. One direction for future study is to build a model with suitable complexity that captures interactions between pathogens and yet is not too detailed to impair inference. Another direction is to refine the error correction to account for the high variability in evaluated likelihoods.

References

- [1] World Health Organization, Special Programme for Research, Training in Tropical Diseases, World Health Organization. Department of Control of Neglected Tropical Diseases, World Health Organization. Epidemic, & Pandemic Alert. (2009). Dengue: guidelines for diagnosis, treatment, prevention and control. World Health Organization.
- [2] Halstead, S. B. (2008). Dengue virus-mosquito interactions. *Annu. Rev. Entomol.*, 53, 273-291.
- [3] Murphy, B. R., & Whitehead, S. S. (2011). Immune Response to Dengue Virus and Prospects for a Vaccine*. *Annual review of immunology*, 29, 587-619.
- [4] Reich, N. G., Shrestha, S., King, A. A., Rohani, P., Lessler, J., Kalayanarooj, S., & Cummings, D. A. (2013). Interactions between serotypes of dengue highlight epidemiological impact of cross-immunity. *Journal of The Royal Society Interface*, 10(86), 20130414.
- [5] Adams, B., Holmes, E. C., Zhang, C., Mammen, M. P., Nimmannitya, S., Kalayanarooj, S., & Boots, M. (2006). Cross-protective immunity can account for the alternating epidemic pattern of dengue virus serotypes circulating in Bangkok. *Proceedings of the National Academy of Sciences*, 103(38), 14234-14239.
- [6] Wearing, H. J., & Rohani, P. (2006). Ecological and immunological determinants of dengue epidemics. *Proceedings of the National Academy of Sciences*, 103(31), 11802-11807.
- [7] Arulampalam, M. S., Maskell, S., Gordon, N., & Clapp, T. (2002). A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *Signal Processing, IEEE Transactions on*, 50(2), 174-188.
- [8] Van Leeuwen, P. J. (2009). Particle filtering in geophysical systems. *Monthly Weather Review*, 137(12), 4089-4114.
- [9] Doucet, A., & Johansen, A. M. (2009). A tutorial on particle filtering and smoothing: Fifteen years later. *Handbook of Nonlinear Filtering*, 12(656-704), 3.
- [10] Kantas, N., Doucet, A., Singh, S. S., Maciejowski, J., & Chopin, N. (2015). On particle methods for parameter estimation in state-space models. *Statistical science*, 30(3), 328-351.
- [11] Andrieu, C., Doucet, A., & Holenstein, R. (2010). Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3), 269-342.
- [12] Ionides, E. L., Bretó, C., & King, A. A. (2006). Inference for nonlinear dynamical systems. *Proceedings of the National Academy of Sciences*, 103(49), 18438-18443.
- [13] Ionides EL, Bhadra A, Atchade Y, King A (2011) Iterated filtering. *Ann Stat* 39:1776–1802.
- [14] Ionides, E. L., Nguyen, D., Atchadé, Y., Stoev, S., & King, A. A. (2015). Inference for dynamic and latent variable models via iterated, perturbed Bayes maps. *Proceedings of the National Academy of Sciences*, 112(3), 719-724.
- [15] Hsiao, K., Miller, J., & de Plinval-Salgues, H. (2005). Particle Filters and Their Applications. *Cognitive Robotics*, April.
- [16] LEE, A. and WHITELEY, N. (2014). Forest resampling for distributed sequential Monte Carlo. Preprint. Available at arXiv:1406.6010.
- [17] King, A. A., Ionides, E. L., Pascual, M., & Bouma, M. J. (2008). Inapparent infections and cholera dynamics. *Nature*, 454(7206), 877-880.
- [18] Lavine, J. S., King, A. A., Andreasen, V., & Bjørnstad, O. N. (2013). Immune boosting explains regime-shifts in prevaccine-era pertussis dynamics. *PloS one*, 8(8), e72086.

- [19] He, D., Ionides, E. L., & King, A. A. (2009). Plug-and-play inference for disease dynamics: measles in large and small populations as a case study. *Journal of the Royal Society Interface*.
- [20] Laneri, K., Bhadra, A., Ionides, E. L., Bouma, M., Dhiman, R. C., Yadav, R. S., & Pascual, M. (2010). Forcing versus feedback: epidemic malaria and monsoon rains in northwest India. *PLoS Comput Biol*, 6(9), e1000898.
- [21] Blake, I. M., Martin, R., Goel, A., Khetsuriani, N., Everts, J., Wolff, C., ... & Grassly, N. C. (2014). The role of older children and adults in wild poliovirus transmission. *Proceedings of the National Academy of Sciences*, 111(29), 10604-10609.
- [22] King, A. A., de Cellès, M. D., Magpantay, F. M., & Rohani, P. (2015, May). Avoidable errors in the modelling of outbreaks of emerging pathogens, with special reference to Ebola. In *Proc. R. Soc. B* (Vol. 282, No. 1806, p. 20150347). The Royal Society.
- [23] Shrestha, S., King, A. A., & Rohani, P. (2011). Statistical inference for multi-pathogen systems. *PLoS Comput Biol*, 7(8), e1002135.
- [24] Shrestha, S., Foxman, B., Weinberger, D. M., Steiner, C., Viboud, C., & Rohani, P. (2013). Identifying the interaction between influenza and pneumococcal pneumonia using incidence data. *Science translational medicine*, 5(191), 191ra84-191ra84.
- [25] King, A. A., Nguyen, D., & Ionides, E. L. (2015). Statistical inference for partially observed Markov processes via the R package pomp. *arXiv preprint arXiv:1509.00503*.
- [26] Ionides, E. L., et al. "Monte Carlo profile confidence intervals for dynamic systems." *Journal of The Royal Society Interface* 14.132 (2017): 20170126.