

# Supplement to “Inference for dynamic and latent variable models via iterated, perturbed Bayes maps”

Edward L. Ionides<sup>1</sup>, D. Nguyen<sup>1</sup>, Y. Atchadé<sup>1</sup>, S. Stoev<sup>1</sup> and A. A. King<sup>2</sup>

July 21, 2015

<sup>1</sup> Department of Statistics and <sup>2</sup>Department of Ecology & Evolutionary Biology,  
The University of Michigan, Ann Arbor, Michigan, USA.

email: ionides@umich.edu, nguyenxd@umich.edu, yvesa@umich.edu, sstoev@umich.edu, kingaa@umich.edu

## Supplementary Content

<b>S1 Weak convergence for occupation measures</b>	<b>S-2</b>
<b>S2 Iterated importance sampling</b>	<b>S-3</b>
<b>S3 Gaussian and near-Gaussian analysis of iterated importance sampling</b>	<b>S-4</b>
<b>S4 A class of exact non-Gaussian limits for iterated importance sampling</b>	<b>S-6</b>
<b>S5 Applying PMCMC to the cholera model</b>	<b>S-7</b>
<b>S6 Applying Liu &amp; West’s method to the toy example</b>	<b>S-8</b>
<b>S7 Consequences of perturbing parameters for the numerical stability of SMC</b>	<b>S-10</b>
<b>S8 Checking conditions B1 and B2</b>	<b>S-10</b>
<b>S9 Additional details for the proof of Theorem 1</b>	<b>S-12</b>
<b>S10 Parameters and parameter ranges for the cholera model</b>	<b>S-16</b>

## S1 Weak convergence for occupation measures

We study the convergence of the processes  $\{W_\sigma(t), 0 \leq t \leq 1\}$  toward  $\{W(t), 0 \leq t \leq 1\}$  as  $\sigma \rightarrow 0$  for Theorem 2. We are interested in showing that the fraction of time  $\{W_\sigma(t)\}$  spends in a set  $\Theta_0 \subset \Theta$  over the discrete set of times  $\{k\sigma^2, k = 1, \dots, 1/\sigma^2\}$  converges in distribution to the fraction of time  $\{W(t)\}$  spends in  $\Theta_0$ . We choose  $\{W_\sigma(t)\}$  to be a right-continuous step function approximation to a diffusion to simplify the relationship between the occupancy fraction over the discrete set of times and over the continuous interval. However, this simplification requires us to work with convergence to  $\{W(t)\}$  in a space of processes with discontinuous sample paths, leading us to work with a Skorokhod topology.

Let  $D_p[0, 1]$  be the space of  $\mathbb{R}^p$ -valued functions on  $[0, 1]$  which are right-continuous with left limits. Let  $X = \{X(t)\}_{t \in [0, 1]}$  and  $\{X_n(t)\}_{t \in [0, 1]}$ ,  $n \geq 1$ , be stochastic processes with paths in  $D_p[0, 1]$ . Let  $\Rightarrow$  denote weak convergence, and suppose that  $X_n \Rightarrow X$  as  $n \rightarrow \infty$  in  $D_p[0, 1]$  equipped with the strong Skorokhod  $J_1$  topology [1].

**Proposition S1** (Proposition VI.1.17 of [1]). *If  $X$  has continuous paths, then  $X_n \Rightarrow X$  as  $n \rightarrow \infty$  in the space  $D_p[0, 1]$  equipped with the uniform metric.*

Suppose that  $f : \mathbb{R}^p \rightarrow \mathbb{R}$  is Borel measurable function and define the map  $T_f : D_p[0, 1] \rightarrow \mathbb{R}$

$$T_f(x) := \int_0^1 f(x(t)) dt, \quad x \in D_p[0, 1].$$

Now, let  $\text{Disc}(T_f)$  denote the set of discontinuity points of  $T_f$ , let  $C_p[0, 1]$  be the space of  $\mathbb{R}^p$ -valued continuous functions on  $[0, 1]$ , and write  $\text{Leb}$  for Lebesgue measure.

**Proposition S2.** *Suppose that  $f$  is bounded. We have that*

$$\text{Disc}(T_f) \cap C_p[0, 1] \subset \left\{ x \in C[0, 1] : \text{Leb}(\{t \in [0, 1] : x(t) \in \text{Disc}(f)\}) > 0 \right\} =: D_f. \quad (\text{S1})$$

*Proof.* Suppose that  $x \in C_p[0, 1]$  does not belong to the right-hand side of (S1) and let  $x_n \rightarrow x$  in  $J_1$ . Then, according to a standard property of the Skorokhod  $J_1$  topology [1] we also have  $\sup_{t \in [0, 1]} |x_n(t) - x(t)| \rightarrow 0$ , as  $n \rightarrow \infty$ . Now, since  $x \notin D_f$ , we have that for almost all  $t \in [0, 1]$ , the point  $x(t)$  is a continuity point of  $f$ . Therefore,  $f(x_n(t)) \rightarrow f(x(t))$ ,  $n \rightarrow \infty$ , for almost all  $t \in [0, 1]$ . Since  $f$  is bounded, the Lebesgue dominated convergence theorem then yields

$$T_f(x_n) \equiv \int_0^1 f(x_n(t)) dt \longrightarrow \int_0^1 f(x(t)) dt \equiv T_f(x), \quad \text{as } n \rightarrow \infty.$$

This completes the proof. □

In the context of stochastic processes, by the Continuous Mapping Theorem, we have convergence in distribution,

$$T_f(X_n) \xrightarrow{d} T_f(X), \quad \text{as } n \rightarrow \infty,$$

provided  $X$  has continuous paths and  $\mathbb{P}(X \in \text{Disc}(f)) = 0$ . In the case when  $f(x) = 1_A(x)$ , the latter translates to

$$\mathbb{P}\{\text{The measure of the time } X \text{ spends on the boundary of } A \text{ is zero}\} = 1. \quad (\text{S2})$$

If the stochastic process has continuous marginal distribution and the set  $A$  has zero boundary, the Fubini's theorem readily implies (S2). Indeed, the probability in (S2) equals

$$\int_{\Omega} \int_0^1 1_{\partial A}(X(t, \omega)) dt \mathbb{P}(d\omega) = \int_0^1 \mathbb{P}(X(t) \in \partial A) dt = 0,$$

provided that  $\text{Leb}(\partial A) = 0$  and if  $X(t)$  has a marginal density for each  $t \in (0, 1)$ . The above arguments lead to the proof of the following result.

**Lemma S1.** *Suppose that  $X_n \Rightarrow X$  in  $D_p[0, 1]$ , equipped with the uniform convergence topology. If the process  $X$  takes values in  $C_p[0, 1]$  and has continuous marginal distributions, then for all bounded Borel functions  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ , that are continuous almost everywhere, i.e. such that  $\text{Leb}(\text{Disc}(f)) = 0$ , we have*

$$\int_0^1 f(X_n(t)) dt \xrightarrow{d} \int_0^1 f(X(t)) dt, \quad \text{as } n \rightarrow \infty.$$

## S2 Iterated importance sampling

When  $N = 1$  in IF2, we obtain a general latent variable algorithm in which each iteration involves importance sampling but not filtering. This situation is called iterated importance sampling [2] and we call this special case of our algorithm IIS2. Iterated importance sampling has previously been used to provide a route into proving convergence of iterated filtering [2, 3]. However, in this article we found it more convenient to prove the full result for iterated filtering directly. Although IIS2 may have some independent value as a practical algorithm, our only use of IIS2 in this article is to provide a convenient environment for explicit computations for Gaussian models in Section S3 and non-Gaussian models in Section S4.

---

### Algorithm IIS2. Iterated importance sampling

---

**input:**

Simulator for $f_X(x; \theta)$	Evaluator for $f_{Y X}(y x; \theta)$
Data, $y^*$	Number of iterations, $M$
Initial parameter swarm, $\{\Theta_j^0, j \text{ in } 1:J\}$	Number of particles, $J$
Perturbation density, $h(\theta \varphi; \sigma)$	Perturbation sequence, $\sigma_{1:M}$

**output:** Final parameter swarm,  $\{\Theta_j^M, j \text{ in } 1:J\}$

For  $m$  in  $1:M$

$\Phi_j^m \sim h(\theta | \Theta_j^{m-1}; \sigma_m)$  for  $j$  in  $1:J$

$X_j^m \sim f_X(x; \Phi_j^m)$  for  $j$  in  $1:J$

$w_j^m = f_{Y|X}(y^* | X_j^m; \Phi_j^m)$  for  $j$  in  $1:J$

Draw  $k_{1:J}$  with  $\mathbb{P}(k_j = i) = w_{n,i}^m / \sum_{u=1}^J w_{n,u}^m$

$\Theta_j^m = \Phi_{k_j}^m$  for  $j$  in  $1:J$

End For

---

A general latent variable model can be specified by a joint density  $f_{XY}(x, y; \theta)$ , with  $X$  taking values in  $\mathbb{X} \subset \mathbb{R}^{\dim(\mathbb{X})}$ ,  $Y$  taking values in  $\mathbb{Y} \subset \mathbb{R}^{\dim(\mathbb{Y})}$  and  $\theta$  taking values in  $\Theta \subset \mathbb{R}^{\dim(\Theta)}$ . The data consist of a single observation,  $y^* \in \mathbb{Y}$ . The likelihood function is

$$\ell(\theta) = f_Y(y^*; \theta) = \int f_{XY}(x, y^*; \theta) dx,$$

and we look for a maximum likelihood estimate (MLE), i.e., a value  $\hat{\theta}$  maximizing  $\ell(\theta)$ . The parameter perturbation step of Algorithm IIS2 is a Monte Carlo approximation to a perturbation map  $H_\sigma$  where

$$H_\sigma g(\theta) = \int g(\varphi) h(\theta | \varphi; \sigma) d\varphi. \quad (\text{S3})$$

A natural choice for  $h(\cdot | \varphi; \sigma)$  is the multivariate normal density with mean  $\varphi$  and variance  $\sigma^2 \Sigma$  for some covariance matrix  $\Sigma$ , but in general  $h$  could be any condition density parameterized by  $\sigma$ . The resampling step of Algorithm IIS2 is a Monte Carlo approximation to a Bayes map,  $B$ , given by

$$Bf(\theta) = f(\theta) \ell(\theta) \left\{ \int f(\varphi) \ell(\varphi) d\varphi \right\}^{-1}. \quad (\text{S4})$$

When the standard deviation of the parameter perturbations is held fixed at  $\sigma_m = \sigma > 0$ , Algorithm IIS2 is a Monte Carlo approximation to  $T_\sigma^M f(\theta)$  where

$$T_\sigma f(\theta) = BH_\sigma f(\theta) = \frac{\int f(\varphi) \ell(\theta) h(\theta | \varphi; \sigma) d\varphi}{\iint f(\varphi) \ell(\xi) h(\xi | \varphi; \sigma) d\varphi d\xi}. \quad (\text{S5})$$

### S3 Gaussian and near-Gaussian analysis of iterated importance sampling

The convergence results of Theorems 1 and 2 in the main text are not precise about the rate of convergence, either toward the MLE as  $\sigma \rightarrow 0$  or toward the stationary distribution as  $M \rightarrow \infty$ . Explicit results are available in the Gaussian case and are also relevant to near-Gaussian situations. The near-Gaussian situation may arise in practice, since the parameter perturbations can be constructed to follow a Gaussian distribution and the log likelihood surface may be approximately quadratic due to asymptotic behavior of the likelihood for large sample sizes. The near-Gaussian situation for a POMP model does not require that the POMP itself is near Gaussian, only that the log likelihood surface is near quadratic. Here, we consider only the univariate case, and only for iterated importance sampling. We offer this simplified case as an illustrative example, rather than an alternative justification for the use of our algorithm. In principle, these results can be generalized, but such results do not add much to the general convergence guarantees already obtained.

We investigate the eigenvalues and eigenfunctions for a Gaussian system, and then we appeal to continuity of the eigenvalues to study systems that are close to Gaussian. Here, we consider the case of a scalar parameter,  $\dim(\Theta) = 1$ , and an additive perturbation given by

$$h(\theta | \varphi; \sigma) = \kappa(\theta - \varphi). \quad (\text{S6})$$

We first study the unnormalized version of (S5) defined as

$$Sf(\theta) = [f(\theta) \ell(\theta)] * \kappa(\theta) = \int [f(\theta - \varphi) \ell(\theta - \varphi)] \kappa(\varphi) d\varphi. \quad (\text{S7})$$

This is a linear map, and we obtain the eigenvalues and eigenfunctions when  $\ell$  and  $h$  are Gaussian in Proposition S3. Iterations of the corresponding normalized map,  $T_\sigma$ , converge to the normalized eigenfunction corresponding to the largest eigenvalue of  $S$ , which can be seen by postponing normalization until having carried out a large number of iterations of the unnormalized map. Suppose, without loss of generality, that the maximum of the likelihood is at  $\theta = 0$ . Let  $\phi(\theta; \sigma)$  be the normal density with mean zero and variance  $\sigma^2$ .

**Proposition S3.** *Let  $S_0$  be the map constructed as in (S7) with the choices  $\ell(\theta) = \phi(\theta; \tau)$  and  $\kappa(\theta) = \phi(\theta; \sigma)$ . Let*

$$u^2 = \left( \sigma^2 + \sqrt{\sigma^4 + 4\sigma^2\tau^2} \right) / 2 = \sigma\tau + o(\sigma). \quad (\text{S8})$$

*The eigenvalues of  $S_0$  are*

$$\lambda_n = \sigma\tau\sqrt{2\pi} \left( \frac{u^2 - \sigma^2}{u^2} \right)^{(n+1)/2},$$

*for  $n = 0, 1, 2, \dots$ , and the corresponding eigenfunctions have the form*

$$e_n = p_n(\theta)\phi(\theta; u), \quad (\text{S9})$$

*where  $p_n$  is a polynomial of degree  $n$ .*

*Proof.* Let  $P_n$  be the subspace of functions of the form  $q(\theta)\phi(\theta; u)$  where  $q$  is a polynomial of degree less than or equal to  $n$ . We show that  $S_0$  maps  $P_n$  into itself, and look at what happens to terms of degree  $n$ . Let  $H_n$  be the Hermite polynomial of degree  $n$ , defined by  $(d/d\theta)^n \phi(\theta; 1) = (-1)^n H_n(\theta)\phi(\theta; 1)$ . Let  $\alpha = (1/u^2 + 1/\tau^2)^{-1/2}$ , and set

$$f(\theta) = \alpha^{-2n} H_n(\theta/\alpha)\phi(\theta; u). \quad (\text{S10})$$

Then,

$$f(\theta)\ell(\theta) = \frac{\alpha}{\sigma\tau\sqrt{2\pi}} \alpha^{-2n} H_n(\theta/\alpha)\phi(\theta; \alpha) = \frac{\alpha}{\sigma\tau\sqrt{2\pi}} (-1)^n \frac{d^n}{d\theta^n} \phi(\theta; \alpha). \quad (\text{S11})$$

Since  $[(d/d\theta)^n f\ell] * \kappa = (d/d\theta)^n [(f\ell) * \kappa]$ , we get

$$(f\ell) * \kappa = \frac{\alpha}{\sigma\tau\sqrt{2\pi}} (-1)^n \frac{d^n}{d\theta^n} \phi(\theta; u) = \frac{\alpha}{\sigma\tau\sqrt{2\pi}} u^{-2n} H_n(\theta/u)\phi(\theta; u). \quad (\text{S12})$$

Writing  $H_n(\theta) = h_0 + h_1\theta + \dots + h_n\theta^n$ , we see that the coefficient of the term in  $\theta^n$  in (S10) is  $\alpha^{-n}h_n$ , whereas in (S12) it is  $\frac{\alpha}{\sigma\tau\sqrt{2\pi}}u^{-n}$ . We have shown that  $S_0$  operating on  $P_n$  multiplies the coefficient of degree  $n$  by a factor of  $\lambda_n$ . Letting  $L_n$  be the matrix representing  $S_0$  on  $P_n$  with the basis  $b_0, \dots, b_m$  given by  $b_m(\theta) = \theta^m\phi(\theta; u)$ , we see that  $L_n$  is lower triangular with diagonal entries  $\lambda_0, \dots, \lambda_n$ . Therefore, the eigenvalues are  $\lambda_0, \dots, \lambda_n$ , and the eigenfunction corresponding to  $\lambda_m$  is in  $P_m$ .  $\square$

The case where  $\log \ell(\theta)$  is close to quadratic is relevant due to asymptotic log quadratic properties of the likelihood function. Choosing  $\kappa(\theta)$  to be Gaussian, as in Proposition S3, we have the following approximation result.

**Proposition S4.** *Let  $S_\epsilon$  be a map as in (S7), with  $\ell$  satisfying  $\sup_\theta |\ell(\theta) - \phi(\theta; \tau)| < \epsilon$  and  $\kappa(\theta) = \phi(\theta; \sigma)$ . For  $\epsilon$  small, the largest eigenvalue of  $S_\epsilon$  is close to  $\lambda_0$  and the corresponding eigenfunction is close to  $\phi(\theta; u)$ .*

*Proof.* Write  $\ell(\theta) = \phi(\theta; \tau) + \eta(\theta)$ , with  $\sup_{\theta} |\eta(\theta)| < \epsilon$ . Then,

$$\|S_{\epsilon}f - S_0f\| = \|(f\eta) * \kappa\| \leq \|f\eta\| \leq \epsilon\|f\|. \quad (\text{S13})$$

Here,  $\|\cdot\|$  is the  $L^2$  norm of a function or the corresponding operator norm (largest absolute eigenvalue). Convolution with  $\kappa$  is a contraction in  $L^2$ , which is apparent by taking Fourier transforms and making use of Parseval's relationship, since all frequencies are shrunk by multiplying with the Fourier transform of  $\kappa$ . From (S13), we have  $\|S_0 - S_{\epsilon}\| < \epsilon$ . This implies that  $S_{\epsilon}$  has a largest eigenvalue  $\mu_0$  with  $|\mu_0 - \lambda_0| < \epsilon$ , based on the representation that

$$|\mu_0| = \|S\| = \sup_f \frac{\|S_{\epsilon}f\|}{\|f\|}. \quad (\text{S14})$$

Writing the corresponding unit eigenfunction as  $w_0$ , we have

$$w_0 = (1/\mu_0)S_{\epsilon}w_0 = (1/\mu_0)[S_0w_0 + \eta], \quad (\text{S15})$$

where  $\|\eta(\theta)\| < \epsilon$ . Writing  $w_0 = \sum_{i=1}^{\infty} \alpha_i e_i$ , in terms of  $\{e_i\}$  from (S9), equation (S15) gives

$$\sum_{i=1}^{\infty} \alpha_i e_i = \sum_{i=1}^{\infty} \alpha_i \frac{\lambda_i}{\mu_0} e_i + \eta = \sum_{i=1}^{\infty} \alpha_i \frac{\lambda_i}{\lambda_0} e_i + \tilde{\eta}, \quad (\text{S16})$$

where  $\|\tilde{\eta}\| < \epsilon(1 + [\lambda_0(\lambda_0 - \epsilon)]^{-1})$ . Comparing terms in  $e_i$ , we see that all terms  $\alpha_1, \alpha_2, \dots$  must be of order  $\epsilon$ .  $\square$

## S4 A class of exact non-Gaussian limits for iterated importance sampling

We look for exact solutions to the equation  $Tf = f$  where  $T = BH$ , as specified in (S5) with  $h(\theta|\varphi; \sigma) = \kappa(\theta - \varphi)$ . This situation corresponds to iterated importance sampling with additive parameter perturbations that have no dependence on  $\sigma$ , as in equation (S6). Now, for  $g(x)$  being a probability density on  $\Theta$ , define

$$\ell_g(x) = c \frac{g(x)}{\kappa * g(x)}, \quad (\text{S17})$$

where  $c$  is a non-negative constant. For likelihood functions of the form (S17), supposing that  $\ell_g$  is integrable, we obtain an eigenfunction  $e(x) = \kappa * g(x)$  for the unnormalized map  $S$  defined in (S7) via the following calculation:

$$Se(x) = c \int \frac{g(x-u)}{(g * \kappa)(x-u)} (g * \kappa)(x-u) \kappa(u) du \quad (\text{S18})$$

$$= c \int g(x-u) \kappa(u) du \quad (\text{S19})$$

$$= c[g * \kappa(x)] = ce(x). \quad (\text{S20})$$

Under conditions such as Theorem 1, it follows that  $\kappa * g$  is the unique eigenfunction for  $T$ , up to a scale factor, and that  $\lim_{M \rightarrow \infty} T^M f = e$ . We do not anticipate practical applications for the conjugacy relationship we have established between the pair  $(\ell_g, \kappa)$  since we see no reason why the likelihood should have the form of (S17). However, this situation does serve to identify a range of possible limiting behaviors for  $T^M$ .

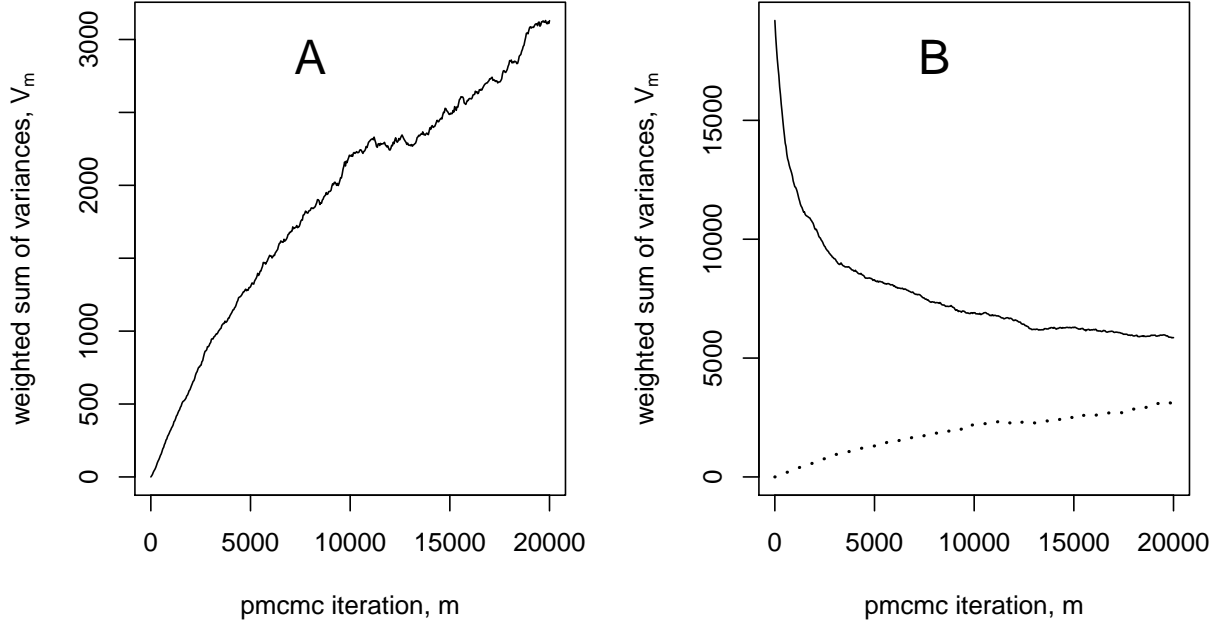


Figure S-1: PMCMC convergence assessment, using the diagnostic quantity in equation S21. (A) Underdispersed chains, all started at the MLE. (B) Overdispersed chains, started with draws from the prior (solid line), and underdispersed chains (dotted line). The average acceptance probability was 0.04238, with Monte Carlo standard error 0.00072, calculated from iterations 5000 through 20000 for the 100 underdispersed PMCMC chains. For the overdispersed chains, the average acceptance probability was 0.04243 with standard error 0.00100.

## S5 Applying PMCMC to the cholera model

We carried out PMCMC for the cholera model, with the prior being uniform on the hyper-rectangle specified by  $\theta_{\text{low}}$  and  $\theta_{\text{high}}$  in Table 1. Thus, the IF1 and IF2 searches were conducted starting with random draws from this prior. Since PMCMC is known to be computationally demanding, we investigated a simplified challenge: investigating the posterior distribution starting at the MLE. This would be appropriate, for example, if one aimed to obtain Bayesian inferences using PMCMC but giving it a helping hand by first finding a good starting value obtained by a maximization procedure. We used the PMMH implementation of PMCMC in `pomp` [4] with parameter proposals following a Gaussian random walk with standard deviations given by  $(\theta_{\text{high}} - \theta_{\text{low}})/100$ . We started 100 independent chains at the estimated MLE in Table 1. Each PMCMC chain, with  $J = 1500$  particles at each of  $M = 2 \times 10^4$  likelihood evaluations, took around 30 hours to run on a single core of the University of Michigan Flux cluster. Writing  $V_{m,d}$  for the sample variance of variable  $d \in \{1, \dots, \dim(\Theta)\}$  among the 100 chains at time  $m \in \{1 \dots, M\}$ , and  $\tau_d$  for the Gaussian random

walk standard deviation for parameter  $d$ , we tracked the quantity

$$V_m = \sum_{d=1}^{\dim(\Theta)} \frac{V_{m,d}}{\tau_d^2}. \quad (\text{S21})$$

Supposing the posterior variance is finite, a necessary requirement for convergence to stationarity as  $m$  increased is for  $V_m$  to approach its asymptotic limit. Since all the chains start at the same place, one expects  $V_m$  to increase toward this limit. The number of iterations required for  $V_m$  to stabilize therefore provides a lower bound on the time taken for convergence of the chain. This test assesses the capability of the chain to explore the region of parameter space with high posterior probability density, rather than the capability to search for this region from a remote starting point. We also tested PMCMC on a harder challenge, investigating convergence of the MCMC chain to its stationary distribution from over-dispersed starting values. We repeated the computation described above, with 100 chains initialized at draws from the prior distribution. The results are shown in Figure S-1. From Figure S-1A, we see that the stationary distribution has not yet been approached for the chains starting at the MLE, since the variance of independent chains continues to increase up to  $M = 2 \times 10^4$ . As a harder test, the variance for the initially overdispersed independent chains should approach that for the initially underdispersed chains, but we see in Figure S-1B that much more computation would be required to achieve this with the algorithmic settings used.

The PMCMC chains used here involved  $JMN = (1.5 \times 10^3) \times (2 \times 10^4) \times (6 \times 10^2) = 1.8 \times 10^{10}$  calls to the dynamic process simulator (the dominating computational expense), and yet failed to converge. By contrast, IF2 with  $JMN = (10^4) \times 10^2 \times (6 \times 10^2) = 6 \times 10^8$  calls to the dynamic process simulator was shown to be an effective tool for global investigation of the likelihood surface. As with all numerical comparisons, it is hard to assess whether poor performance is a consequence of poor algorithmic choices. Conceptually, a major difference between iterated filtering and PMCMC is that the filtering particles in IF2 investigate the parameter space and latent dynamic variable space simultaneously, whereas in PMCMC each filtering iteration is used only to provide a single noisy likelihood evaluation. It may not be surprising that algorithms such as PMCMC struggle in situations where filtering is a substantial computational expense and the likelihood surface is sufficiently complex that many thousands of Monte Carlo steps are required to explore it. Indeed, IF1 and IF2 remain the only algorithms that have currently been demonstrated computationally capable of efficient likelihood-based inference for situations of comparable difficulty to our example.

## S6 Applying Liu & West’s method to the toy example

Bayesian parameter estimation for POMP models using sequential Monte Carlo with perturbed parameters was proposed by [6]. Similar approaches using alternative nonlinear filters have also been widely used [7, 8]. Liu & West [5] proposed a development on the approach of [6] which combines parameter perturbations with a contraction that is designed to counterbalance the variation added by the perturbations, thereby approximating the posterior distribution of the parameters for the fixed parameter model of interest. Liu & West [5] also included an auxiliary particle filter procedure in their algorithm [9]. The auxiliary particle filter is a version of sequential Monte Carlo which looks ahead to a future observation when deciding which particles to propagate. Generally, auxiliary particle filter algorithms do not have the plug-and-play property [10, 11] since they involve constructing weights that require evaluation of the transition density for the latent process.



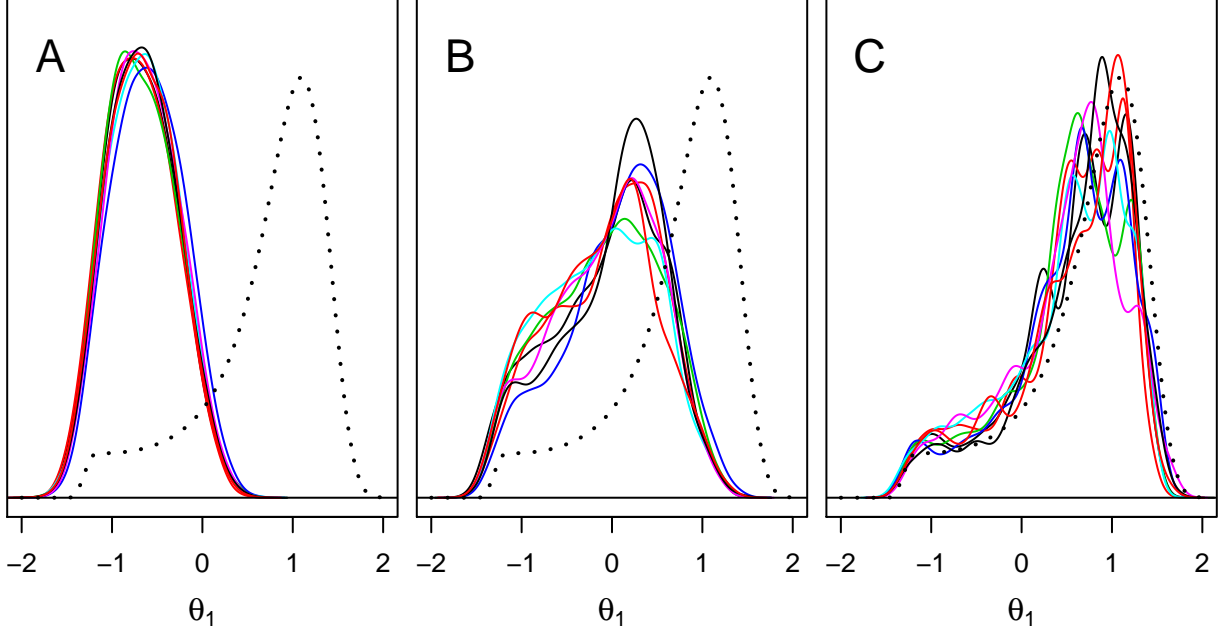


Figure S-2: The Liu & West algorithm [5] applied to the toy example with varying values of the discount factor: (A)  $\delta = 0.99$ ; (B)  $\delta = 0.999$ ; (C)  $\delta = 0.9999$ . Solid lines show 8 independent estimates of the marginal posterior density of  $\theta_1$ . The black dotted line shows the true posterior density.

In addition, the auxiliary particle filter does not necessarily have superior performance over a basic sequential Monte Carlo filter [12]. To compare with IF2 and PMCMC on our toy example, we therefore employ a version of the Liu & West algorithm, which we call LW, that omits the auxiliary particle filter procedure. LW carries out the key innovation of parameter perturbation and contraction (Steps 3 and 4 in Sec. 10.4 of [5]) while omitting the auxiliary particle filter (Steps 1 and 2, and the denominator in Step 5, in Sec. 10.4 of [5]). LW was implemented via the `bsmc2` function of the `pomp` package [4]. If an effective auxiliary particle filter were available for a specific computation, it could also be used to enhance other sequential Monte Carlo based inference procedures such as IF1, IF2 and PMCMC.

For the numerical results reported in Fig. S-2 we used  $J = 10^4$  particles for LW. This awards the same computational resources to LW that we gave IF1 and IF2 for the results in Fig. 1. The magnitude of the perturbations in LW is controlled by a discount factor ( $\delta$  in the notation of [5]), and we considered three values,  $\delta \in \{0.99, 0.999, 0.9999\}$ . Liu & West [5] suggested that  $\delta$  should take values in the range  $\delta \in [0.95, 0.99]$ , with smaller values of  $\delta$  reducing Monte Carlo variability while increasing bias in the approximation to the target posterior distribution. For our toy example, we see from Fig. S-2A that the choice  $\delta = 0.99$  results in a stable Monte Carlo computation (since all eight realizations are close). However, Fig. S-2A also reveals a large amount of bias. Increasing  $\delta$  to 0.999, Fig. S-2B shows some increase in the Monte Carlo variability and some decrease in the bias. Further increasing  $\delta$  to 0.9999, Fig. S-2C shows the bias becomes small while the Monte Carlo variability continues to increase. Values of  $\delta$  very close to one are numerically tractable for this toy model, but not in most applications. As  $\delta$  approaches one, the ensuing numerical instability exemplifies the principal reason why Bayesian and likelihood-based inference for POMP models is

challenging despite the development of modern nonlinear filtering techniques.

The justification provided by [5] for their algorithm is based on a Gaussian approximation to the posterior distribution. Specifically, [5] argued that the posterior distribution should be approximately unchanged by carrying out a linear contraction toward its mean followed by adding an appropriate perturbation. Therefore, it may be unsurprising that LW performs poorly in the presence of nonlinear ridges in the likelihood surface. Other authors have reported poor numerical performance for the algorithm of [5], e.g., Fig. 2 of [13] and Fig. 2 of [14]. Our results are consistent with these findings, and we conclude that the approach of [5] should be used with considerable caution when the posterior distribution is not close to Gaussian.

## S7 Consequences of perturbing parameters for the numerical stability of SMC

The IF2 algorithm applies sequential Monte Carlo (SMC) to an extended POMP model in which the time-varying parameters are treated as dynamic state variables. This procedure increases the dimension of the state space by the number of time-varying parameters. Empirically, SMC has been found effective in many low dimensional systems but its numerical performance can degrade in larger systems. A natural concern, therefore, is the extent to which the extension of the state variable in IF2 increases the numerical challenge of carrying out SMC effectively. Two rival heuristics suggest different answers. One intuitive (but not universally correct) argument is that adding variability to the system stabilizes numerically unstable filtering problems, since it gives each particle at least a slim chance of following a trajectory compatible with the data. An opposing intuition, that SMC breaks down rapidly as the dimension increases, has theoretical support [15]. However, the theoretical arguments of [15] may be driven more by increasing the observation dimension than increasing the state dimension, so their relevance in the present situation is not entirely clear.

We investigated numerical stability of SMC, in the context of our cholera example, by measuring the effective sample size (ESS) [16]. We investigated the ESS for two parameter vectors, the MLE and an alternative value for which SMC is more numerically challenging. We carried out particle filtering with and without random walk perturbations to the parameters, obtaining the results presented in Fig. S-3. We found that the random walk perturbations led to a 5% decrease in the average ESS at the MLE, but a 13% increase in the average ESS at the alternative parameter vector. This example demonstrates that the random walk perturbations can have both a cost and a benefit for numerical stability, with the benefit outweighing the cost as the filtering problem becomes more challenging.

## S8 Checking conditions B1 and B2

We check B1 and B2 when  $\Theta$  is a rectangular region in  $\mathbb{R}^{\dim(\Theta)}$ , with  $h_n(\theta|\phi;\sigma)$  describing a Gaussian random walk having as a limit a reflected Brownian motion on  $\Theta$ . A more general study of the limit of reflected random walks to reflected Brownian motions (in particular, including limits where the random walk step distribution satisfies B5) was presented by Bossy et al. [17]. The specific examples of the IF2 algorithm given in our paper all employ Gaussian random walk perturbations for the parameters. The examples did not employ boundary conditions to constrain the parameter to a bounded set. While such conditions could be used to ensure practical stability

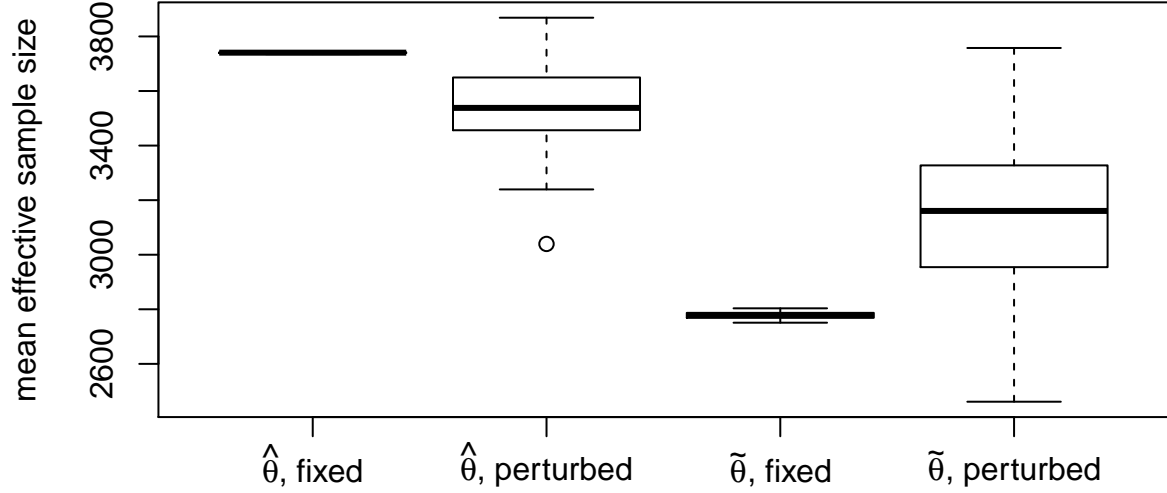


Figure S-3: Effective sample size (ESS) for SMC with fixed parameters and with perturbed parameters. We ran SMC for the cholera model with the parameter vector set at the MLE,  $\hat{\theta}$ , and at an alternative parameter vector  $\tilde{\theta}$  for which the first 18 parameters in Table 1 were multiplied by a factor of 0.8. We defined the ESS at each time point by the reciprocal of the sum of squares of the normalized weights of the particles. The mean ESS was calculated as the average of these ESS values over the 600 time points. Repeating this computation 100 times, using  $J = 10^4$  particles, gave 100 mean ESS values shown in the “fixed” columns of the box-and-whisker plot. Repeating the computation with additional parameter perturbations having random walk standard deviation of 0.01 gave the 100 mean ESS values shown in the “perturbed” column. For both parameter vectors, the perturbations greatly increase the spread of the mean ESS. At  $\hat{\theta}$ , the perturbations decreased the mean ESS value by 5% on average, whereas at  $\tilde{\theta}$  the perturbations increased the mean ESS value by 13% on average. The MLE may be expected to be a favorable parameter value for stable filtering, and our interpretation is that the parameter perturbations have some chance of moving the SMC particles away from this favorable region. When started away from the MLE, the numerical stability of the IF2 algorithm benefits from the converse effect that the parameter perturbations will move the SMC particles preferentially toward this favorable region. For parameter values even further from the MLE than  $\tilde{\theta}$ , SMC may fail numerically for a fixed parameter value yet be feasible with perturbed parameters.

of the algorithm, we view the conditions primarily as a theoretical device to assist the mathematical analysis of the algorithm.

Suppose that  $\Theta = [a_1, b_1] \times [a_2, b_2] \times \dots \times [a_{\dim(\Theta)}, b_{\dim(\Theta)}]$ . For each coordinate direction  $d = 1, \dots, \dim(\Theta)$ , let  $R_d : \mathbb{R} \rightarrow [a_d, b_d]$  be the reflection map defined recursively by

$$R_d(x) = \begin{cases} x & \text{if } x \in [a_d, b_d] \\ R_d(2b_d - x) & \text{if } x > b_d \\ R_d(2a_d - x) & \text{if } x < a_d \end{cases}.$$

Let  $h_{n,d}(\theta_d | \phi_d; \sigma)$  be the density of  $R_d(\phi_d + \sigma Z)$  where  $Z$  is a standard Normal random variable. Let  $h_n(\theta | \phi; \sigma)$  be the joint density corresponding to the product of  $h_{n,1}, \dots, h_{n,\dim(\Theta)}$ . This choice of  $h_n$  corresponds to a perturbation process for the parameter vector in the IF2 algorithm following a Gaussian random walk on  $\Theta$  with reflective boundary conditions, independently in each coordinate direction. By construction, the finite dimensional distributions of  $W_\sigma(t)$  at the set of times

$$\{k\sigma^2 : k = 0, 1, 2, \dots \text{ and } k\sigma^2 \leq 1\}$$

exactly match the corresponding finite dimensional distributions of a reflected Brownian motion  $\{W(t)\}$  taking values in  $\Theta$ . This  $\{W(t)\}$  gives a construction of the limiting process whose existence is assumed in B1. For  $A \subset \Theta$ , we see from this construction of  $\{W(t)\}$  that the probability  $\{W(t)\}$  is in  $A$  for all  $\epsilon \leq t \leq 1$  is greater than the corresponding probability for an unreflected Brownian motion,  $\{W_{(u)}(t)\}$  with the same intensity parameter. It is routine to check that  $\{W_{(u)}(t)\}$  has a positive probability of remaining in any open set  $A$  for all  $\epsilon \leq t \leq 1$  uniformly over all values of  $W_{(u)}(0) \in \Theta$ . Thus, we have completed the check of condition B1.

To check B2, the positivity of the marginal density of  $W(t)$  on  $\Theta$ , uniformly over the value of  $W(0)$ , again follows since this density is larger than the known density for  $W_{(u)}(t)$ .

## S9 Additional details for the proof of Theorem 1

In the main text, a condensed proof of Theorem 1 is provided to describe the key steps in the argument. Here, we restate Theorem 1 and provide a more detailed proof. The reader is referred back to the main text for the notation and statement of conditions B2 and B4.

**Theorem 1.** *Let  $T_\sigma$  be the map defined by [1] in the main text, and suppose B2 and B4. There exists a unique probability density  $f_\sigma$  such that for any probability density  $f$  on  $\Theta$ ,*

$$\lim_{m \rightarrow \infty} \|T_\sigma^m f - f_\sigma\|_1 = 0, \quad (\text{S22})$$

where  $\|f\|_1$  is the  $L^1$  norm of  $f$ . Let  $\{\Theta_j^M, j = 1, \dots, J\}$  be the output of IF2, with  $\sigma_m = \sigma > 0$ . There exists a finite constant  $C$  such that

$$\limsup_{M \rightarrow \infty} \mathbb{E} \left[ \left| \frac{1}{J} \sum_{j=1}^J \phi(\Theta_j^M) - \int \phi(\theta) f_\sigma(\theta) d\theta \right| \right] \leq \frac{C \sup_\theta |\phi(\theta)|}{\sqrt{J}}. \quad (\text{S23})$$

*Proof.* Let  $L^1(\Theta)$  denote the space of integrable real-valued functions on  $\Theta$  with norm  $\|f\|_1 = \int |f(\theta)| d\theta$ . For non-negative measures  $\mu$  and  $\nu$  on  $\Theta$ , let  $\|\mu - \nu\|_{\text{tv}}$  denote the total variation distance and let  $H(\mu, \nu)$  denote the Hilbert metric distance [18, 19]. The measures  $\mu$  and  $\nu$  are said to be comparable if they are both nonzero and there exist constants  $0 < a \leq b$  such that  $a\nu(A) \leq \mu(A) \leq b\nu(A)$  for all measurable subsets  $A \subset \Theta$ . For comparable measures,  $H(\mu, \nu)$  is defined by

$$H(\mu, \nu) = \log \frac{\sup_A \mu(A)/\nu(A)}{\inf_A \mu(A)/\nu(A)}, \quad (\text{S24})$$

with the supremum and infimum taken over measurable subsets  $A \subset \Theta$  having  $\nu(A) > 0$ . For noncomparable measures, the Hilbert metric is defined by  $H(0, 0) = 0$  and otherwise  $H(\mu, \nu) = \infty$ . The Hilbert metric is invariant to multiplication by a positive scalar,  $H(a\mu, \nu) = H(\mu, \nu)$ . This projective property makes the Hilbert metric convenient to investigate the Bayes map: in the context of the following proof, the projective property lets us analyze the linear map  $S_\sigma$  to study the nonlinear map  $T_\sigma$ .

For  $\theta_{0:N} \in \Theta^{N+1}$ , we single out the last component of  $\theta_{0:N}$  by writing  $\check{\ell}(\theta_{0:N}) = \check{\ell}(\theta_{0:N-1}, \theta_N)$  and  $h(\theta_{0:N} | \phi) = h(\theta_{0:N-1}, \theta_N | \phi)$ . Then, for  $\phi$  and  $\theta$  in  $\Theta$ , we define

$$s_\sigma(\phi, \theta) = \int h(\theta_{0:N-1}, \theta | \phi, \sigma) \check{\ell}(\theta_{0:N-1}, \theta) d\theta_{0:N-1}. \quad (\text{S25})$$

The function  $s_\sigma$  in (S25) defines a linear operator  $S_\sigma f(\theta) = \int s_\sigma(\phi, \theta) f(\phi) d\phi$  that maps  $L^1(\Theta)$  into itself. Notice that  $T_\sigma f(\theta) = S_\sigma f(\theta) / \|S_\sigma f\|_1$ . More generally, if  $\mu$  is a probability measure on  $\Theta$ ,  $S_\sigma \mu$  denotes the function  $S_\sigma \mu(\theta) = \int s_\sigma(\phi, \theta) \mu(d\phi)$ . Notice also that  $S_\sigma^m f$ , the  $m$ -th iterate of  $S_\sigma$ , can be written as  $S_\sigma^m f(\theta) = \int s_\sigma^{(m)}(\phi, \theta) f(\phi) d\phi$ , where  $s_\sigma^{(1)}(\phi, \theta) = s_\sigma(\phi, \theta)$ , and for  $m \geq 2$ ,  $s_\sigma^{(m)}(\phi, \theta) = \int s_\sigma(\phi, u) s_\sigma^{(m-1)}(u, \theta) du$ . Using the definition of  $\check{\ell}$  and B4,

$$\begin{aligned} s_\sigma(\phi, \theta) &= \int h(\theta_{0:N-1}, \theta | \phi, \sigma) \int f_X(x_{0:N} | \theta_{0:N-1}, \theta) f_{Y|X}(y_{1:N}^* | x_{0:N}) dx_{0:N} d\theta_{0:N-1} \\ &\geq \epsilon^N \int h(\theta_{0:N-1}, \theta | \phi, \sigma) d\theta_{0:N-1}, \end{aligned} \quad (\text{S26})$$

and, similarly,

$$s_\sigma(\phi, \theta) \leq \epsilon^{-N} \int h(\theta_{0:N-1}, \theta | \phi, \sigma) d\theta_{0:N-1}. \quad (\text{S27})$$

By iterating the inequalities (S26) and (S27), assumption B2 implies that there exists  $m_0 \geq 1$  such that for any  $m \geq m_0$ , there exist  $0 < \delta_m < \infty$ , a probability measure  $\lambda_m$  on  $\Theta$  such that for all measurable subsets  $A \subset \Theta$  and all  $\theta \in \Theta$ ,

$$\delta_m \lambda_m(A) \leq \int_A s^{(m)}(\theta, \phi) d\phi \leq \delta_m^{-1} \lambda_m(A). \quad (\text{S28})$$

In other words,  $S_\sigma^{m_0}$  is mixing in the sense of [19]. In the terminology of [18], this means that for each  $m \geq m_0$ ,  $S_\sigma^m$  has finite projective diameter (see Lemma 2.6.2 of [18]). Therefore, by Theorem 2.5.1 of [18], we conclude that  $S_\sigma$  has a unique non-negative eigenfunction  $f_\sigma$  with  $\|f_\sigma\|_1 = 1$ , and for any density  $f$  on  $\Theta$ , as  $q \rightarrow \infty$ ,

$$\left\| \frac{[S_\sigma^{m_0}]^q f}{\|[S_\sigma^{m_0}]^q f\|_1} - f_\sigma \right\|_1 = \|T_\sigma^{m_0 q} f - f_\sigma\|_1 \rightarrow 0.$$

This implies the statement (S22), by writing for any  $m \geq 1$ ,  $m = qm_0 + r$ , for  $0 \leq r \leq m_0 - 1$ , and  $T_\sigma^m f = [T_\sigma^{qm_0} T_\sigma^r f]$ .

Let the initial particle swarm  $\{\Theta_j^0, 1 \leq j \leq J\}$  consist of independent draws from the density  $f$ . To prove (S23), we decompose  $M = qm_0 + r$ , for some  $r \in \{0, \dots, m_0 - 1\}$ , and we introduce the empirical measures  $\mu^{(0)} = J^{-1} \sum_{j=1}^J \delta_{\Theta_j^{(r)}}$ , and for  $k = 1, \dots, q$ ,  $\mu^{(k)} = J^{-1} \sum_{j=1}^J \delta_{\Theta_j^{(r+m_0k)}}$ , so that  $\mu^{(q)} = J^{-1} \sum_{j=1}^J \delta_{\Theta_j^{(M)}}$ . We then write, for any bounded measurable function  $\phi$ ,

$$\begin{aligned} \mu^{(q)}(\phi) - [T_\sigma^M f](\phi) &= \mu^{(q)}(\phi) - [T_\sigma^{m_0q} \mu^{(0)}](\phi) + [T_\sigma^{m_0q} \mu^{(0)}](\phi) - [T_\sigma^{m_0q} T_\sigma^r f](\phi) \\ &= \sum_{i=1}^q \left\{ [T_\sigma^{m_0(i-1)} \mu^{(q-i+1)}](\phi) - [T_\sigma^{m_0i} \mu^{(q-i)}](\phi) \right\} \\ &\quad + [T_\sigma^{m_0q} \mu^{(0)}](\phi) - [T_\sigma^{m_0q} T_\sigma^r f](\phi). \end{aligned}$$

Using Theorem 2 of [20], we can find a finite constant  $C_3$  such that for all  $k \geq 1$ , and writing  $\|\phi\|_\infty = \sup_\theta |\phi(\theta)|$ ,

$$\rho = \sup_{\phi: \|\phi\|_\infty=1} \mathbb{E} \left[ \left| \mu^{(k)}(\phi) - [T_\sigma^{m_0k} \mu^{(k-1)}](\phi) \right| \right] \leq \frac{C_3}{\sqrt{J}}, \quad (\text{S29})$$

with B4 implying that the constant  $C_3$  constructed by [20] does not depend on  $\mu^{(k-1)}$ . Since  $S_\sigma^{m_0}$  is mixing and (S28) holds, using Lemma 3.4, Lemma 3.5, Lemma 3.8 and Equation (7) of [19], we have

$$\begin{aligned} \mathbb{E} \left[ \left| [T_\sigma^{m_0q} \mu^{(0)}](\phi) - [T_\sigma^{m_0q} T_\sigma^r f](\phi) \right| \right] &\leq \|\phi\|_\infty \mathbb{E} \left[ \|T_\sigma^{m_0q} \mu^{(0)} - T_\sigma^{m_0q} T_\sigma^r f\|_{\text{tv}} \right] \\ &\leq \frac{2\|\phi\|_\infty}{\log 3} \mathbb{E} \left[ H \left( S_\sigma^{m_0q} \mu^{(0)}, S_\sigma^{m_0q} T_\sigma^r f \right) \right] \\ &\leq \frac{2\|\phi\|_\infty}{\log 3} \left( \frac{1 - \delta_{m_0}^2}{1 + \delta_{m_0}^2} \right)^{q-2} \frac{1}{\delta_{m_0}^2} \mathbb{E} \left[ \|T_\sigma^{m_0} \mu^{(0)} - T_\sigma^{m_0} T_\sigma^r f\|_{\text{tv}} \right] \\ &\leq \frac{4\|\phi\|_\infty}{\log 3} \left( \frac{1 - \delta_{m_0}^2}{1 + \delta_{m_0}^2} \right)^{q-2} \frac{1}{\delta_{m_0}^2} \frac{\rho}{\delta_{m_0}^2}. \end{aligned}$$

For  $i = 3, \dots, q$ , a similar calculation gives

$$\begin{aligned} \mathbb{E} \left[ \left| T_\sigma^{m_0(i-1)} \mu^{(q-i+1)}(\phi) - T_\sigma^{m_0i} \mu^{(q-i)}(\phi) \right| \right] &= \mathbb{E} \left[ \left| T_\sigma^{m_0(i-1)} \mu^{(q-i+1)}(\phi) - T_\sigma^{m_0(i-1)} T_\sigma^{m_0} \mu^{(q-i)}(\phi) \right| \right] \\ &\leq \frac{4\|\phi\|_\infty}{\log 3} \left( \frac{1 - \delta_{m_0}^2}{1 + \delta_{m_0}^2} \right)^{i-3} \frac{1}{\delta_{m_0}^2} \frac{\rho}{\delta_{m_0}^2}. \end{aligned}$$

The case  $i = 1$  boils down to (S29), where the case  $i = 2$  gives by similar calculations:

$$\mathbb{E} \left[ \left| T_\sigma^{m_0} \mu^{(q-1)}(\phi) - T_\sigma^{2m_0} \mu^{(q-2)}(\phi) \right| \right] \leq 2\|\phi\|_\infty \frac{\rho}{\delta_{m_0}^2}.$$

Hence, using (S29),

$$\mathbb{E} \left[ \left| \mu^{(q)}(\phi) - [T_\sigma^M f](\phi) \right| \right] \leq \frac{C_3 \|\phi\|_\infty}{\sqrt{J}} \left( 1 + \frac{2}{\delta_{m_0}^2} + \frac{4}{\log 3} \left( \frac{1}{\delta_{m_0}^2} \right)^{2^{q-2}} \sum_{j=0}^{2^{q-2}} \left( \frac{1 - \delta_{m_0}^2}{1 + \delta_{m_0}^2} \right)^j \right).$$

We conclude that there exists a finite constant  $C_4$  such that

$$\mathbb{E} \left[ \left| \frac{1}{J} \sum_{j=1}^J \phi(\Theta_j^M) - \int \phi(\theta) [T_\sigma^M f](\theta) d\theta \right| \right] \leq \frac{C_4 \|\phi\|_\infty}{\sqrt{J}}. \quad (\text{S30})$$

Equation (S23) follows by combining (S30) with (S22). □

# S10 Parameters and parameter ranges for the cholera model

Table S-1. Parameters for the cholera model.

	$\hat{\theta}$	$\theta_{\text{low}}$	$\theta_{\text{high}}$
$\gamma$	20.80	10.00	40.00
$\epsilon$	19.10	0.20	30.00
$m$	0.06	0.03	0.60
$\beta_{\text{trend}} \times 10^2$	-0.50	-1.00	0.00
$\beta_1$	0.75	-4.00	4.00
$\beta_2$	6.38	0.00	8.00
$\beta_3$	-3.44	-4.00	4.00
$\beta_4$	4.23	0.00	8.00
$\beta_5$	3.33	0.00	8.00
$\beta_6$	4.55	0.00	8.00
$\omega_1$	-1.69	-10.00	0.00
$\omega_2$	-2.54	-10.00	0.00
$\omega_3$	-2.84	-10.00	0.00
$\omega_4$	-4.69	-10.00	0.00
$\omega_5$	-8.48	-10.00	0.00
$\omega_6$	-4.39	-10.00	0.00
$\sigma$	3.13	1.00	5.00
$\tau$	0.23	0.10	0.50
$S_0$	0.62	0.00	1.00
$I_0$	0.38	0.00	1.00
$R_{1,0}$	0.00	0.00	1.00
$R_{2,0}$	0.00	0.00	1.00
$R_{3,0}$	0.00	0.00	1.00

$\hat{\theta}$  is the MLE reported by [21]. Three parameters were fixed ( $\delta = 0.02$ ,  $N_s = 6$  and  $k = 3$ ) following [21]. Units are  $\text{year}^{-1}$  for  $\gamma$ ,  $\epsilon$ ,  $m$ ,  $\beta_{\text{trend}}$  and  $\delta$ ; all other parameters are dimensionless.  $\theta_{\text{low}}$  and  $\theta_{\text{high}}$  are the lower and upper bounds for a hyper-rectangle used to generate starting points for the search. Non-negative parameters ( $\gamma$ ,  $\epsilon$ ,  $m$ ,  $\sigma$ ,  $\tau$ ) were logarithmically transformed for optimization. Unit scale parameters ( $S_0$ ,  $I_0$ ,  $R_{1,0}$ ,  $R_{2,0}$ ,  $R_{3,0}$ ) were optimized on a logistic scale. These parameters were rescaled using the known population size to give the initial state variables, e.g.,  $S(t_0) = S_0\{S_0 + I_0 + R_{1,0} + R_{2,0} + R_{3,0}\}^{-1}P(t_0)$ .



## Supplementary References

- [1] Jacod, J & Shiryaev, A. N. (1987) *Limit theorems for stochastic processes*. (Springer-Verlag, Berlin).
- [2] Ionides, E. L, Bhadra, A, Atchadé, Y, & King, A. A. (2011) Iterated filtering. *Annals of Statistics* **39**, 1776–1802.
- [3] Doucet, A, Jacob, P. E, & Rubenthaler, S. (2013) Derivative-free estimation of the score vector and observed information matrix with application to state-space models. *Arxiv*, <http://arxiv.org/abs/1304.5768>.
- [4] King, A. A, Ionides, E. L, Bretó, C. M, Ellner, S, & Kendall, B. (2009) pomp: Statistical inference for partially observed markov processes. *R package*, available at <http://cran.r-project.org/web/packages/pomp>.
- [5] Liu, J & West, M. (2001) in *Sequential Monte Carlo Methods in Practice*, eds. Doucet, A, de Freitas, N, & Gordon, N. J. (Springer, New York), pp. 197–224.
- [6] Kitagawa, G. (1998) A self-organising state-space model. *Journal of the American Statistical Association* **93**, 1203–1215.
- [7] Anderson, B. D & Moore, J. B. (1979) *Optimal Filtering*. (Prentice-Hall, New Jersey).
- [8] Wan, E & van der Merwe, R. (2000) *The unscented Kalman filter for nonlinear estimation*. (IEEE), pp. 153–158.
- [9] Pitt, M. K & Shepard, N. (1999) Filtering via simulation: Auxillary particle filters. *Journal of the American Statistical Association* **94**, 590–599.
- [10] Bretó, C, He, D, Ionides, E. L, & King, A. A. (2009) Time series analysis via mechanistic models. *Annals of Applied Statistics* **3**, 319–348.
- [11] He, D, Ionides, E. L, & King, A. A. (2010) Plug-and-play inference for disease dynamics: Measles in large and small towns as a case study. *Journal of the Royal Society Interface* **7**, 271–283.
- [12] Johansen, A. M & Doucet, A. (2008) A note on the auxiliary particle filter. *Statistics and Probability Letters* **78**, 1498–1504.
- [13] Storvik, G. (2002) Particle filters for state-space models with the presence of unknown static parameters. *IEEE Transactions on Signal Processing* **50**, 281–289.
- [14] Chopin, N, Jacob, P. E, & Papaspiliopoulos, O. (2013) SMC<sup>2</sup>: an efficient algorithm for sequential analysis of state space models. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **75**, 397–426.
- [15] Bengtsson, T, Bickel, P, & Li, B. (2008) in *Probability and Statistics: Essays in Honor of David A. Freedman*, eds. Speed, T & Nolan, D. (Institute of Mathematical Statistics, Beachwood, OH), pp. 316–334.

- [16] Liu, J. S. (2001) *Monte Carlo Strategies in Scientific Computing*. (Springer, New York).
- [17] Bossy, M, Gobet, E, & Talay, D. (2004) A symmetrized Euler scheme for an efficient approximation of reflected diffusions. *Journal of Applied Probability* pp. 877–889.
- [18] Eveson, S. P. (1995) Hilbert’s projective metric and the spectral properties of positive linear operators. *Proceedings of the London Mathematical Society* **3**, 411–440.
- [19] Le Gland, F & Oudjane, N. (2004) Stability and uniform approximation of nonlinear filters using the Hilbert metric and application to particle filters. *Annals of Applied Probability* **14**, 144–187.
- [20] Crisan, D & Doucet, A. (2002) A survey of convergence results on particle filtering methods for practitioners. *IEEE Transactions on Signal Processing* **50**, 736–746.
- [21] King, A. A, Ionides, E. L, Pascual, M, & Bouma, M. J. (2008) Inapparent infections and cholera dynamics. *Nature* **454**, 877–880.