



Journal of the American Statistical Association

Publication details, including instructions for authors and subscription information:

<http://www.tandfonline.com/loi/uasa20>

Malaria in Northwest India: Data Analysis via Partially Observed Stochastic Differential Equation Models Driven by Lévy Noise

Anindya Bhadra, Edward L. Ionides, Karina Laneri, Mercedes Pascual, Menno Bouma and Ramesh C. Dhiman

A. Bhadra is Doctoral Student and E. L. Ionides is Associate Professor, Department of Statistics, University of Michigan, Ann Arbor, MI 48109. K. Laneri is Postdoctoral Student, Department of Ecology & Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109. M. Pascual is Rosemary Grant Collegiate Professor, Department of Ecology & Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109 and Investigator, Howard Hughes Medical Institute, 4000 Jones Bridge Road, Chevy Chase, MD 20815-6789. M. Bouma is Adjunct Faculty, Department of Infectious and Tropical Diseases, London School of Hygiene and Tropical Medicine, University of London, London, United Kingdom. R. C. Dhiman is Research Scientist, National Institute of Malaria Research, Sector 8, Dwarka, Delhi-110077, India. This work was funded in part by support from the RAPIDD program at the Science & Technology Directorate of the Department of Homeland Security and the Fogarty International Center, National Institutes of Health; the National Science Foundation (DMS-0805533); the National Oceanic and Atmospheric Administration (Oceans and Health Program NA 04O AR 460019); and the Graham Environmental Sustainability Institute of the University of Michigan. The authors thank the editor (Hal Stern), the associate editor, and two anonymous referees for their helpful comments.

Version of record first published: 24 Jan 2012.

To cite this article: Anindya Bhadra, Edward L. Ionides, Karina Laneri, Mercedes Pascual, Menno Bouma and Ramesh C. Dhiman (2011): Malaria in Northwest India: Data Analysis via Partially Observed Stochastic Differential Equation Models Driven by Lévy Noise, *Journal of the American Statistical Association*, 106:494, 440-451

To link to this article: <http://dx.doi.org/10.1198/jasa.2011.ap10323>

PLEASE SCROLL DOWN FOR ARTICLE

Full terms and conditions of use: <http://www.tandfonline.com/page/terms-and-conditions>

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden.

The publisher does not give any warranty express or implied or make any representation that the contents will be complete or accurate or up to date. The accuracy of any instructions, formulae, and drug doses should be independently verified with primary sources. The publisher shall not be liable for any loss, actions, claims, proceedings, demand, or costs or damages whatsoever or howsoever caused arising directly or indirectly in connection with or arising out of the use of this material.

Malaria in Northwest India: Data Analysis via Partially Observed Stochastic Differential Equation Models Driven by Lévy Noise

Anindya BHADRA, Edward L. IONIDES, Karina LANERI, Mercedes PASCUAL, Menno BOUMA, and Ramesh C. DHIMAN

Many biological systems are appropriately described by partially observed Markov process (POMP) models, also known as state space models. Such models also arise throughout the physical and social sciences, in engineering, and in finance. Statistical challenges arise in carrying out inference on nonlinear, nonstationary, vector-valued POMP models. Methodologies that depend on the Markov process model only through numerical solution of sample paths are said to have the plug-and-play property. This property enables consideration of models for which the evaluation of transition densities is problematic. Our case study employs plug-and-play methodology to investigate malaria transmission in Northwest India. We address the scientific question of the respective roles of environmental factors, immunity, and nonlinear disease transmission dynamics in epidemic malaria. Previous debates on this question have been hindered by the lack of a statistical investigation that gives simultaneous consideration to the roles of human immunity and the fluctuations in mosquito abundance associated with environmental or ecological covariates. We present the first time series analysis integrating these various components into a single vector-valued dynamic model. We are led to investigate a POMP involving a system of stochastic differential equations driven by Lévy noise. We find a clear role for rainfall and evidence to support models featuring the possibility of clinical immunity.

An online supplement presents details of the methodology implemented and two additional figures.

KEY WORDS: Iterated filtering; Partially observed Markov process; *Plasmodium falciparum*; Sequential Monte Carlo.

1. INTRODUCTION

Malaria is currently a widespread tropical and subtropical disease, with approximately 500 million cases per year (Snow et al. 2005) resulting in over one million deaths (Hay et al. 2005). Malaria is caused by infection with a protozoan parasite which is transmitted between humans by mosquitoes. The disease was eliminated from North America and Europe during the first half of the 20th century, primarily by sanitary and agricultural developments which reduced contact between humans and mosquitoes below the level required to sustain disease transmission (Packard 2007). From 1955 to 1969 the World Health Organization ran an ambitious Global Malaria Eradication Program, based on mosquito control by extensive spraying with the insecticide DDT and treatment with the antimalarial drug chloroquine (Packard 2007). In India, malaria incidence declined dramatically during the Global Malaria Eradication Program. A crippling burden of approximately 75 million cases per year was reduced to a reported incidence of 49,151 in 1961 (Kumar et al. 2007). However, rather than continuing this decline, malaria incidence crept up through the 1960s.

The reemergence in India has been attributed to the increasing cost and decreasing supply of DDT, resistance developed by mosquitoes to DDT, and increasing resistance of malaria parasites to chloroquine (Sharma 1996; Kumar et al. 2007). After increasing to over six million reported cases annually in the 1970s, malaria incidence has since stabilized at around two million cases per year (Kumar et al. 2007). These official statistics are an indication of the trend of incidence but fail to include many cases which are treated outside the public health system. A more accurate estimate of recent incidence may be 11 million cases per year (World Health Organization 2008).

Hopes for a global eradication of malaria have recently been raised once more. Eradication has been stated as an explicit goal of the Bill and Melinda Gates Foundation, with the endorsement of the World Health Organization and the Roll Back Malaria Partnership (Roberts and Enserink 2007). The main technologies underpinning this aspiration are long-lasting insecticide-treated bed nets and a new generation of artemisinin-derived antimalarial drugs. Although global eradication is probably unrealistic with currently available tools (Greenwood 2009), there is great potential to reduce the heavy global burden of malaria. One of the lessons learned from the previous eradication program is that effective control requires adaptation to local patterns of disease transmission (Greenwood 2009). Improved quantitative understanding of transmission is therefore a necessary component of control and prevention efforts.

The early mathematical models of Ross (1911) and MacDonald (1957) have long been a foundation for developing malaria control strategies (McKenzie and Samba 2004). Many extensions have been proposed to these mathematical models, allowing for biological aspects such as genetic diversity of the

A. Bhadra is Doctoral Student and E. L. Ionides is Associate Professor (E-mail: ionides@umich.edu), Department of Statistics, University of Michigan, Ann Arbor, MI 48109. K. Laneri is Postdoctoral Student, Department of Ecology & Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109. M. Pascual is Rosemary Grant Collegiate Professor, Department of Ecology & Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109 and Investigator, Howard Hughes Medical Institute, 4000 Jones Bridge Road, Chevy Chase, MD 20815-6789. M. Bouma is Adjunct Faculty, Department of Infectious and Tropical Diseases, London School of Hygiene and Tropical Medicine, University of London, London, United Kingdom. R. C. Dhiman is Research Scientist, National Institute of Malaria Research, Sector 8, Dwarka, Delhi-110077, India. This work was funded in part by support from the RAPIDD program at the Science & Technology Directorate of the Department of Homeland Security and the Fogarty International Center, National Institutes of Health; the National Science Foundation (DMS-0805533); the National Oceanic and Atmospheric Administration (Oceans and Health Program NA 040 AR 460019); and the Graham Environmental Sustainability Institute of the University of Michigan. The authors thank the editor (Hal Stern), the associate editor, and two anonymous referees for their helpful comments.

parasite (Gupta et al. 1994), the mosquito and parasite lifecycle (McKenzie and Bossert 2005), the development of drug resistance (Koella and Antia 2003; Klein et al. 2008), and exposure-dependent partial immunity (Dietz, Molineaux, and Thomas 1974; Aron and May 1982; Filipe et al. 2007). Given the size of the public health issue and the extent of the research into malaria transmission, it may be surprising how few studies investigate the relationship between these dynamic models and available population-level time series data. Investigations relating disease models (which are typically partially observed nonlinear Markov processes) to time series data have a long tradition of inspiring developments in statistical analysis of stochastic dynamic systems (Bartlett 1960; Ellner et al. 1998; Finkenstädt and Grenfell 2000; Ionides, Bretó, and King 2006; Cauchemez et al. 2008). Indeed, the most convenient disease systems to study, such as measles, are still considered a challenge for statistical inference (Cauchemez and Ferguson 2008; He, Ionides, and King 2010). Analysis of measles dynamics is simplified by clear clinical diagnosis, direct human-to-human transmission, lifelong immunity following infection, and the availability of extensive spatio-temporal incidence data. The study of malaria dynamics is hindered by nonspecific symptoms; one usually has to work under the assumption that malaria is the cause of sickness for patients who have a high fever and are found, by inspection of a blood slide under a microscope, to be infected with *Plasmodium* parasites. However, asymptomatic *Plasmodium* infections are not unusual, and there are many alternative potential causes of fever. Second, human immunity to malaria wanes with time and gives varying levels of protection to diverse disease strains. Clinical immunity (i.e., protection to symptomatic infection) can result from repeated infections, and leads to infections with a reduced transmissibility. Third, malaria transmission is dependent on mosquito abundance. Malaria transmission is highly sensitive to the density, longevity and biting habits of the mosquito vector. These entomological quantities vary considerably in space and time, both within and between vector species (Packard 2007). Time series of vector abundance and behavior directly relevant to long-term population-level studies are therefore generally unavailable.

In Section 2, we develop a quantitative approach to relate malaria transmission to available time series data. We aim to construct statistical models of the population-level transmission dynamics which are at once sophisticated enough to capture the important features of the biological system and simple enough that they can be rigorously assessed using available data. Mathematically, our models are a set of coupled nonlinear system of stochastic differential equations driven by Lévy noise. Whereas certain specific models could be constructed using the more usual choice of Gaussian noise, a general framework which satisfies necessary nonnegativity constraints can more readily be built using nonnegative noise built from nondecreasing Lévy processes such as the Gamma process. Lévy process models have been proposed for a range of applications, ranging from option pricing in finance to quantum mechanics (Applebaum 2004). However, statistically efficient inference from general classes of nonstationary partially observed systems driven by Lévy noise has not previously, to our knowledge, been demonstrated. Here, we use the term *statistically efficient* in an informal sense, to describe methodology leading to parameter estimates whose uncertainty approximates that of Bayesian or

likelihood-based estimates. Statistical efficiency becomes an important consideration when building models whose complexity is at, or close to, the limit which the available data can support.

Numerical solution of SDEs driven by general Lévy noise is available using methods extending the well-studied special case of Gaussian noise (Protter and Talay 1997; Jacod 2004). Recently, statistical methodologies for partially observed Markov processes have been proposed for which the dynamic model enters into the inference procedure only through the availability of numerical solutions (i.e., simulated sample paths). Such methodologies are said to have the *plug-and-play* property, since simulation code can be plugged directly the inference algorithm (Bretó et al. 2009; He, Ionides, and King 2010). One might hope that such techniques facilitate statistical inference for models of malaria. Our case study demonstrate that this is indeed the case, by carrying out inference as a routine application of a recently developed likelihood-based plug-and-play technique called iterated filtering (Ionides, Bretó, and King 2006). Other plug-and-play techniques have been proposed in the context of simulated moment methods (McFadden 1989; Kendall et al. 2005; Wood 2010), approximate Bayesian methods (Liu and West 2001; Sisson, Fan, and Tanaka 2007; McKinley, Cook, and Deardon 2009; Wilkinson 2011) and asymptotically exact Bayesian inference (Andrieu, Doucet, and Holenstein 2010). By comparison, standard expectation-maximization and Markov chain Monte Carlo approaches (Cappé, Moulines, and Rydén 2005) require the evaluation of transition densities—which can lead to difficulties, or even complete failure, on continuous time POMP models (Roberts and Stramer 2001). Therefore, the development of plug-and-play methodology promises to greatly extend the classes of dynamic models available for use in data analysis.

Moment-based and approximate Bayesian methods sacrifice some statistical efficiency as a trade-off for the plug-and-play property. Beyond the loss of statistical efficiency, these approximate methods also suffer from a potential lack of objectivity in the choice of approximation. Moreover, methods based on simulated moments have a substantial practical limitation that they cannot be routinely extended to nonstationary models. The obstacle for statistically efficient plug-and-play methods is computational efficiency. Indeed, Wood (2010) has recently argued that statistically efficient inference (both Bayesian and non-Bayesian) is computationally intractable for highly nonlinear partially observed stochastic population models such as our malaria model. Demonstration of the computational feasibility of likelihood-based analysis via iterated filtering provides a counter-example. Existing exact Bayesian plug-and-play methods are computationally expensive compared to iterated filtering (Bhadra 2010) though whether or not alternative plug-and-play techniques could feasibly have been used is outside the scope of this case study.

Section 3 presents a data analysis, through which we aim both to demonstrate our statistical approach and to draw conclusions about the respective roles of immunity and climate variability for epidemic malaria transmission. Epidemic or ‘unstable’ malaria (Kiszewski and Teklehaimanot 2004) occurs when conditions are only occasionally favorable for disease transmission, for example, due to cold or dry seasons which preclude

mosquito activity. Waning of immunity during the absence of exposure to malaria can lead to high levels of severe infection in epidemics. By contrast, the repeated exposures in regions of endemic or 'stable' malaria result in acquisition of immunity that protects from severe forms of the disease. We focus on two questions. First, what is the appropriate degree of model complexity which is necessary to understand population dynamics of epidemic malaria? This issue is basic to developing scientifically acceptable models for malaria which quantitatively match population-level incidence data. Second, what is the role of climate fluctuations, such as interannual changes in rainfall patterns, for determining the interannual variability of disease incidence? Despite agreement on the sensitivity of the mosquito vector to environmental conditions, there has been considerable controversy on the respective roles of environmental forcing versus epidemiological considerations, fueled by the lack of a quantitative statistical approach which can make a formal comparison of rival hypotheses. In particular, for malaria in East African highlands, some investigators have found that interannual variability in rainfall and temperature can explain a substantial share of the variability in regional malaria incidence time series (Pascual et al. 2006), whereas others have proposed that oscillating levels of immunity in the population act as the major driver (Hay et al. 2002). We broaden this specific debate by analyzing data from another unstable malaria transmission environment, in an arid region of Northwest India, where the role of rainfall variability is less controversial but has not been addressed together with immunity in the context of the population dynamics of the disease. It is in desert and highland regions, at the edge of the distribution of the disease, that we expect climate variability and climate change to be potentially most relevant to disease dynamics due to the limiting roles of rainfall and temperature. The data analysis in this article focuses on a newly available malaria incidence time series for the Kutch district, an arid region in the state of Gujarat. The scientific argument is expanded on elsewhere (Laneri et al. 2010), and our primary goal here is to describe the statistical foundations for building and analyzing dynamic models of population-level malaria transmission that can be confronted to time series data.

2. MALARIA TRANSMISSION: A STATISTICAL MODEL

We start by describing some relevant biology; for a more complete introduction we recommend Warrell and Gilles (2002). The unicellular protozoan parasites of the genus *Plasmodium* which cause malaria are transmitted between humans by the female of certain species of *Anopheles* mosquito. The *Plasmodium* lifecycle consists of multiple stages in both human and mosquito hosts. When a mosquito takes a blood meal from an infected human, male and female *Plasmodium* gametocytes may be ingested. Sexual reproduction of the parasite takes place within a vector mosquito's stomach, resulting in the formation of sporozoites which migrate to the mosquito's salivary glands. Upon a subsequent blood meal, the sporozoites can infect another human. They enter the bloodstream, become sequestered in the liver, reemerge into the blood, reproduce asexually in erythrocyte stages, and eventually produce gametocytes to complete the cycle. During the stages in a human host, the *Plasmodium* must do battle with the complex human immune system

which attacks sporozoite, erythrocyte and gametocyte stages (Artavanis-Tsakonas, Tongren, and Riley 2003). The effectiveness of the immune response depends, among other things, on system memory from previous exposure to related parasites. Transmission of malaria relies upon the availability of infected humans, susceptible humans, and mosquitoes having sufficient longevity. The mosquito longevity is critical for the viability of the *Plasmodium* lifecycle since the time taken for the *Plasmodium* to undergo ingestion, reproduction, development and retransmission to a human host is comparable to the mean lifespan of the mosquito.

The majority of severe and fatal human malaria cases are caused by infection with *P. falciparum*. The other widespread species is *P. vivax*, which is characterized by less severe symptoms with the possibility of relapse many months after infection. To develop a quantitative representation, we will write down a model for falciparum malaria (i.e., disease resulting from infection with *P. falciparum*) which captures some key aspects of the human, parasite and vector dynamics. This model could be extended to vivax malaria by the inclusion of relapse. Our goal is to present a *statistical model* in the sense that it is sufficiently parsimonious that the parameters can be estimated directly from available data, as carried out in Section 3.

We divide humans into five distinct classes: S_1 , fully susceptible to infection; S_2 , protected from severe infection, but susceptible to mild reinfection; E , exposed (i.e., carrying *Plasmodium* parasites which have not yet matured into gametocytes); I_1 , infected and gametocytemic; I_2 , possessing a mild, asymptomatic infection with reduced gametocyte levels (Filipe et al. 2007; Klein et al. 2008). An innovative feature of this framework, compared to other epidemiological models previously fitted to population-level time series data, is the inclusion of an explicit representation of the vector dynamics: A mosquito stage λ represents the latent force of infection, capturing the likelihood of successful transmission from human to human together with a distributed delay. This formulation avoids explicit consideration of mosquito abundance, survival and behavior. In other words, we limit our inclusion of vector dynamics to the aspect that is most directly relevant to the human disease.

Figure 1 represents diagrammatically the modeled flows between these classes, formally defined by Equations (1)–(7) below. We write μ_{XY} for the rate of transition from class X to class Y , for X and Y in $\{S_1, S_2, E, I_1, I_2\}$. In addition, we introduce a per-capita birth rate, μ_{BS_1} , into the completely susceptible class. Deaths occur at a constant rate $\mu_{XD} = \delta$ from each class $X \in \{S_1, S_2, E, I_1, I_2\}$. As mortality from acute malarial infection has become small in India, we do not include disease-induced mortality in our model. The total population size $P(t)$ is supposed known by interpolation from the decennial census. Transition from S_2 to I_2 can be interpreted as reinfection with clinical immunity, that is, reduced symptoms which do not lead the patient to seek medical attention (Artavanis-Tsakonas, Tongren, and Riley 2003). We suppose that $\mu_{S_2 I_2} = c \mu_{S_1 E}$ with some constant of proportionality $0 \leq c \leq 1$. Our model also includes the possibility of failing to acquire any protective immunity following infection, by transitioning directly from I_1 back to S_1 without passing through I_2 and S_2 . This can arise through prompt treatment with antimalarial drugs, in which case the body does not have time to build an immune response (Klein

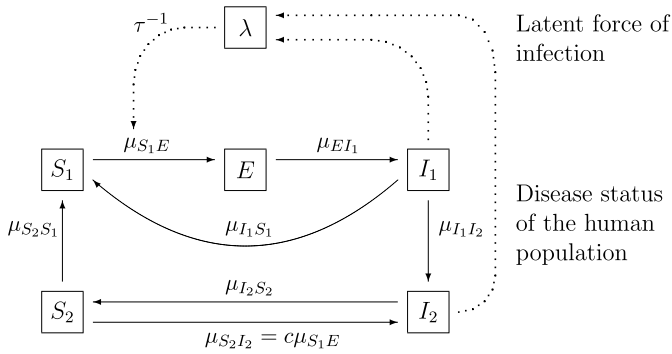


Figure 1. A compartment model of malaria transmission. Human classes are S_1 (susceptible), S_2 (partially protected), E (exposed, carrying a latent infection), I_1 (infected and infectious), and I_2 (asymptomatic, with reduced infectivity). A solid arrow from X to Y denotes the possibility of transition, with rate μ_{XY} . Dotted arrows represent interactions between the human and mosquito stages of the parasite. Mosquito dynamics are modeled via the latent force of infection, λ , which affects the current force of infection μ_{S_1E} with a mean latency time of τ . We call this model VS²EI² with ‘V’ for ‘vector’ followed by the human classes and their multiplicities. In the subcase with $\mu_{I_2S_2} = \infty$ and $\mu_{S_2I_2} = \mu_{I_1S_1} = 0$, the class I_2 can be eliminated to allow direct transitions from I_1 to S_2 , and individuals in S_2 are fully protected. The classes $\{S_1, E, I_1, S_2\}$ can then be mapped onto the classes $\{S, E, I, R\}$ in a standard epidemiological susceptible-exposed-infected-recovered model (Anderson and May 1991) with added vector dynamics and waning immunity, so we call this special case VSEIR.

et al. 2008). Alternatively, it can be a consequence of the necessity for multiple infections before the body learns to mount an effective general-purpose defense against clinical symptoms in the face of the genetic diversity of the *Plasmodium*.

$$\begin{aligned} dS_1/dt = & \mu_{BS_1}P - \mu_{S_1E}(t)S_1 + \mu_{I_1S_1}I_1 \\ & + \mu_{S_2S_1}S_2 - \mu_{S_1D}S_1, \end{aligned} \quad (1)$$

$$dS_2/dt = \mu_{I_2S_2}I_2 - \mu_{S_2S_1}S_2 - \mu_{S_2I_2}S_2 - \mu_{S_2D}S_2, \quad (2)$$

$$dE/dt = \mu_{S_1E}S_1 - \mu_{EI_1}E - \mu_{ED}E, \quad (3)$$

$$dI_1/dt = \mu_{EI_1}E - \mu_{I_1S_1}I_1 - \mu_{I_1I_2}I_1 - \mu_{I_1D}I_1, \quad (4)$$

$$dI_2/dt = \mu_{I_1I_2}I_1 + \mu_{S_2I_2}S_2 - \mu_{I_2S_2}I_2 - \mu_{I_2D}I_2, \quad (5)$$

$$\lambda(t) = \frac{I_1(t) + qI_2(t)}{P(t)} \bar{\beta} \exp\left\{Z_t\beta + \sum_{i=1}^{n_s} \beta_i s_i(t)\right\}, \quad (6)$$

$$\begin{aligned} \mu_{S_1E}(t) = & \int_{-\infty}^t \gamma(t-s)\lambda(s) d\Gamma(s) \\ \text{for } \gamma(t) = & \frac{(k/\tau)^k t^{k-1}}{(k-1)!} \exp\{-kt/\tau\}. \end{aligned} \quad (7)$$

The malarial status of the human population in (1)–(5) corresponds to a large population limit of homogeneous individual-level interactions (Anderson and May 1991; Keeling and Rohani 2008). The integral Equation (7) combines *Plasmodium* development and mosquito survival in the classic Ross–Macdonald model (Macdonald 1957; Aron and May 1982). The gamma-distributed latency, with mean τ and variance τ^2/k , was chosen to allow a differential representation that facilitates numerical solution (Lloyd 2001). Specifically, we define

$\lambda_1(t), \dots, \lambda_k(t)$ to satisfy

$$d\lambda_1/dt = (\lambda - \lambda_1)k\tau^{-1} d\Gamma/dt, \quad (8)$$

$$d\lambda_i/dt = (\lambda_{i-1} - \lambda_i)k\tau^{-1} \quad \text{for } i = 2, \dots, k. \quad (9)$$

Setting $\mu_{S_1E}(t) = \lambda_k(t)$, (8)–(9) is equivalent to (7). Stochasticity in this system is presumed to arise principally from variations in vector abundance and behavior, which is modeled by the stochastic process $\Gamma(t)$. We take $\Gamma(t)$ to be a gamma process representing integrated noise with intensity σ^2 . This is defined as a process with stationary independent increments such that $\Gamma(t) - \Gamma(s) \sim \text{Gamma}([t-s]/\sigma^2, \sigma^2)$ where $\text{Gamma}(a, b)$ is the gamma distribution with mean ab and variance ab^2 . Although $\Gamma(t)$ is a jump process, and therefore its sample paths are not differentiable, one can interpret the process $d\Gamma/dt$ as multiplicative gamma noise (Bretó et al. 2009). The reason to choose gamma noise over the more familiar Gaussian noise is to enforce the positivity of $\mu_{S_1E}(t)$ and hence all the state variables in (1)–(7). Supposing that $d\Gamma(t)/dt$ represents nonnegative white noise is equivalent to requiring that $\Gamma(t)$ be a Lévy jump process with a nonnegative jump distribution (Applebaum 2004). The gamma process was selected for being a relatively simple and well-studied nonnegative Lévy process. The Equations (1)–(6) and (8)–(9) then define a set of coupled stochastic differential equations driven by Lévy noise. We solve this system numerically via the Euler method (Protter and Talay 1997; Jacod 2004) with a time-step of one day. Whereas all state variables in the unavailable exact solution are nonnegative, it is possible for the Euler method to generate numerical approximations violating this constraint. We monitored the frequency of these occurrences; they were rare to the point of negligibility in our analysis.

Equation (6) is based on mass-action principles (Anderson and May 1991; Keeling and Rohani 2008). Here, q represents the transmissibility, relative to full-blown infections, from asymptomatic infections in partially immune individuals. The seasonality of disease transmission is modeled by the coefficients $\{\beta_i\}$ corresponding to a periodic cubic B-spline basis $\{s_i(t), i = 1, \dots, n_s\}$ constructed using n_s evenly spaced knots. Time-varying covariates enter via the row vector Z_t with coefficients in a column vector β . The dimensional constant $\bar{\beta}$ is required to give $\mu_{S_1E}(t)$ units of t^{-1} , and we set $\bar{\beta} = 1 \text{ yr}^{-1}$.

At first inspection, our model may appear to be a dauntingly complex specification based on many assumptions that one cannot hope to validate. However, this work builds on a long history of developing and using similar models (Anderson and May 1991; Keeling and Rohani 2008). All the parameters in (1)–(7) have interpretable scientific meaning and can therefore be discussed in the context of the literature on malaria transmission. Indeed, our model can also be criticized as an oversimplification, since we do not incorporate many of the biological aspects developed in previous models (such as Gupta et al. 1994; Koella and Antia 2003; McKenzie and Bossert 2005; Chitnis, Cushing, and Hyman 2006). In addition, our model does not make allowances for spatial, socioeconomic, age-related and genetic inhomogeneities among the population. Such structure could play an important role. Nevertheless, models based on homogeneous populations are often sufficient to describe the major features of disease transmission dynamics

(Earn et al. 2000). In the face of biological complexity, a major part of the value of constructing and analyzing dynamic models is to develop an understanding of the key components driving the behavior of the biological system. In our modeling framework, alternative model specifications can readily be analyzed and compared, building on the results reported here.

A measurement model provides a formal connection between the dynamic process model and available data. Here we give an abstract representation, deferring concrete discussion of data to Section 3. We write $\{t_n, n = 1, \dots, N\}$ for the times of the N observations, and we suppose that the model is initialized at some time $t_0 < t_1$. We define the number of new cases in the n th interval to be $C_n = \int_{t_{n-1}}^{t_n} \mu_{EI} E(s) ds$. The reported number of confirmed cases, y_n , is then modeled conditional on C_n as $y_n | C_n \sim \text{Negbin}(\rho C_n, \psi^2)$, where $\text{Negbin}(\alpha, \beta)$ is the negative binomial distribution with mean α and variance $\alpha + \alpha^2 \beta$. This distribution allows for the possibility of over-reporting or under-reporting, and can be viewed as an over-dispersed Poisson distribution with dispersion parameter ψ . We refer to ρ as the reporting rate. It is known that only a small fraction of malaria cases are treated in the public clinics which contribute to district statistics (Kumar et al. 2007), so we expect $\rho \ll 1$. The exact interpretation of ρ is necessarily sensitive to the severity of disease that is required to be classified as a case.

Although environmental covariates affect many biological systems, quantifying their dynamic role can be a formidable task, both from a scientific and a statistical perspective (Bjornstad and Grenfell 2001). The flexibility of plug-and-play statistical methodology permits scientific considerations to determine ways in which covariates might appropriately be included in the analysis. Here, we take Z_t to be a scalar covariate measuring the thresholded rainfall integrated over a time interval $[t - u, t]$. Specifically, from the accumulated rainfall data $\{r_n, n = 1, \dots, N\}$ at times t_1, \dots, t_N we interpolated a continuous-time cubic spline $r(t)$ and then set $\tilde{Z}_t = \max\{\int_{t-u}^t r(s) ds - v, 0\}$. This specification is designed to

represent parsimoniously the threshold and lag effects which are to be expected in biological systems. The covariate was standardized by setting $Z_t = (\tilde{Z}_t - \bar{Z})/\sigma_Z$, where $\bar{Z} = (t_N - t_0)^{-1} \int_{t_0}^{t_N} \tilde{Z}_s ds$ and $\sigma_Z^2 = (t_N - t_0)^{-1} \int_{t_0}^{t_N} (\tilde{Z}_s - \bar{Z})^2 ds$. This standardization makes the coefficient β a dimensionless quantity which is expected to vary on a unit scale.

3. DATA ANALYSIS

Figure 2 shows a plot of the monthly confirmed cases of *P. falciparum* and monthly rainfall in the district of Kutch in the state of Gujarat in Northwest India between January 1987 and December 2006. The record of the malaria cases was obtained from the National Institute of Malaria Research in India, and was originally compiled by the office of the District Malaria Officer. The rainfall time series was obtained from a local district weather station run by the Indian Meteorology Department. Rainfall in Kutch is concentrated within the monsoon season, and Kutch experiences the seasonal epidemic malaria typical of arid regions of India (Swaroop 1949; Bouma and van der Kaay 1994; Kiszewski and Teklehaimanot 2004). Visually, a lag relationship, with rainfall leading malaria, may seem evident from this figure. Since rainfall typically peaks during the summer monsoon and malaria typically peaks a few months later, in late fall, one might see the appearance of a lag relationship in the absence of a direct link. The correlation between total monsoon rainfall (aggregated over June–August) and total fall cases (aggregated over October–December) is 0.84 over these 20 years, which is suggestive of a causal relationship. However, the pattern of monsoon rainfall has cycles of 2–4 years which matches cycles that are predicted in malaria due to the building up of population immunity in epidemics followed by subsequent waning of immunity and birth of newly susceptible children (Pascual et al. 2008). This confounding of intrinsic cycles (e.g., immunity) with the effect of extrinsic cycles (e.g., climate variability) adds difficulty to the interpretation of such correlations. Modeling both intrinsic and extrinsic effects simultaneously provides a way to strengthen scientific conclusions. This

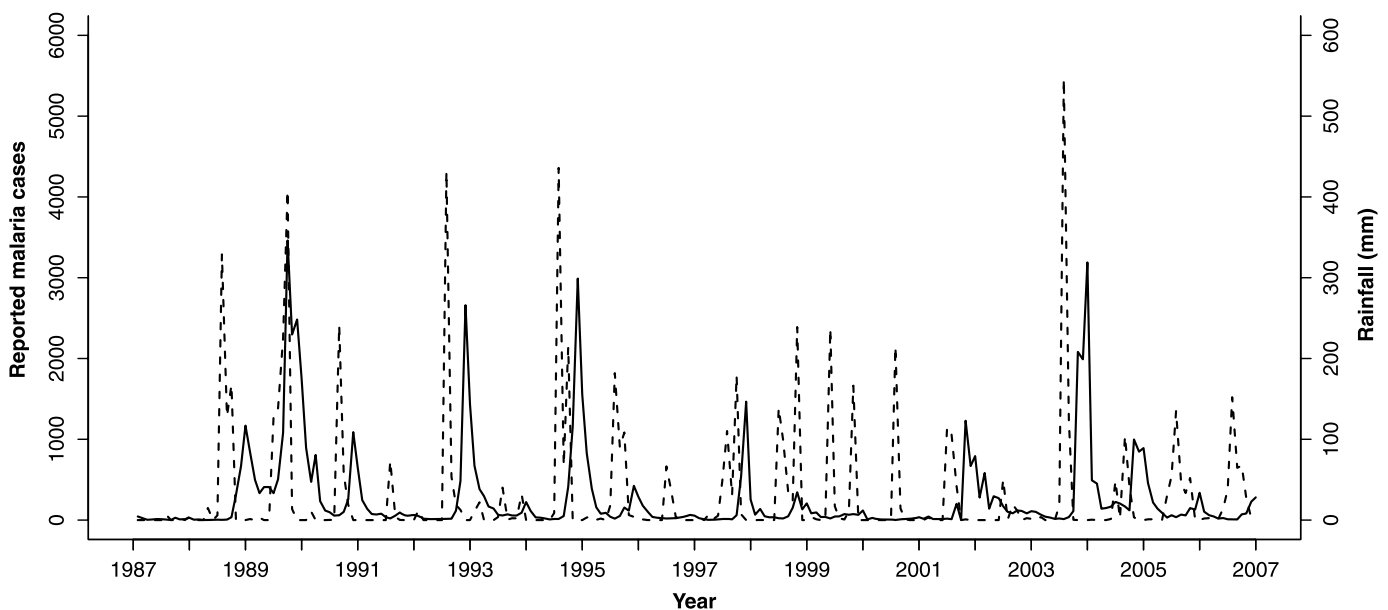


Figure 2. Monthly reported *P. falciparum* malaria cases (solid line) and monthly rainfall from a local weather station (broken line) for Kutch.

is analogous to using multiple regression to control for potential confounding variables, but here we must take into account the nonlinear stochastic feedbacks and lagged relationships in the dynamic system. In this investigation, we fixed the rainfall covariate by constructing \tilde{Z}_t using $u = 5$ mo and $v = 200$ mm. Additional analysis of the relationship with rainfall are published elsewhere (Laneri et al. 2010), but it should be clear that the approach we develop here has the flexibility to address alternative hypotheses concerning this as well as many other questions about malaria dynamics.

We carried out likelihood-based inference via iterated filtering, a plug-and-play sequential Monte Carlo procedure for calculating maximum likelihood estimates which was introduced by Ionides, Bretó, and King (2006). Heuristically, iterated filtering algorithms carry out sequential Monte Carlo while adding stochastic perturbations to the parameters. In subsequent iterations, the intensity of these stochastic perturbations is decreased and so the likelihood surface is investigated at increasingly local scales. We refer the reader to the online supplement, and to the relevant literature (King et al. 2008; Bretó et al. 2009; He, Ionides, and King 2010), for further discussion of iterated filtering methodology. Computer code to generate an Euler solution to the dynamic model and to evaluate the density of the measurement model is all that the user must supply to embark on statistical analysis via general-purpose software implementing such a plug-and-play procedure. Iterated filtering was implemented using the `pomp` software package (King et al. 2009) which encodes the algorithm presented by King et al. (2008, supplementary text). There are tuning parameters which affect the numerical efficiency of the maximization algorithm; the values we used are reported in the online supplement. These algorithmic parameters are inconsequential for the inferential conclusions once numerical convergence has been confirmed by checking consistency over a range of starting values for the likelihood maximization. If all the model parameters share a unit scale of variability, selection of the algorithmic parameters is simplified.

With this in mind, we worked with the logarithmic transform of nonnegative parameters and the logit transform of parameters taking values in the interval $(0, 1)$. On this common scale, standard values of the algorithmic parameters gave acceptable optimization performance. All reported results are transformed back to the original scale.

It is a substantial computational challenge to investigate a multimodal likelihood function, with 25 parameters some combinations of which are weakly identifiable, based on Monte Carlo estimates of the likelihood which involves integrating over all the unobserved state variables at $(t_N - t_0)/\Delta = 20 \times 365$ time points. Iterated filtering, being a Monte Carlo optimization technique, has inherent stochasticity which enables the algorithm to escape from some local maxima. However, to investigate global optimization (both to search for distinct modes of the likelihood function, and to confirm that this search has been adequately carried out) we used a large number of starting values selected randomly from a defined region of the parameter space. A complete description of the procedure, including the region sampled for starting values, is given in the online supplement. Results from this search are presented in Figure 3. The highest value of the likelihood values in Figure 3A is within 0.1 log units of the maximum likelihood obtained during all our searches, including the computationally intensive exercise of computing profile likelihoods for each parameter in the construction of the standard errors in Table 3. The standard error of our sequential Monte Carlo likelihood evaluations is also approximately 0.1 log units. We infer that, despite multimodality ridges in the likelihood surface, a systematic procedure for reliable maximization of the likelihood function using iterated filtering is computationally feasible.

Figure 3B displays one example of multimodality and weak identifiability. The convergence points for μ_{EI} and τ are tightly clustered in two separated groups, suggesting bimodality. These groups have very similar values of $\mu_{EI}^{-1} + \tau$, the total latent period of the *Plasmodium* in human and mosquito hosts combined. When analyzing case report data, it might be expected

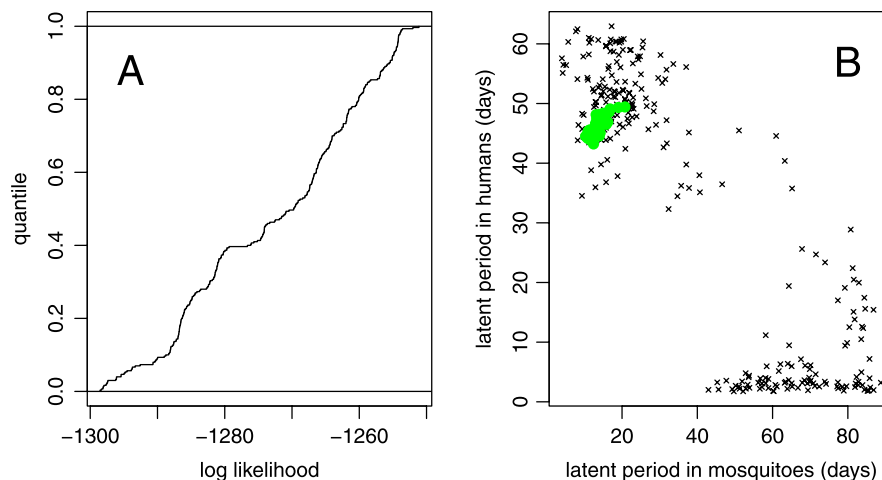


Figure 3. (A) Empirical distribution function of log likelihoods for 300 parameter estimates using independently selected starting values. These results were obtained for Kutch, using the VS^2EI^2 model with rainfall. A complete description of the search procedure is given in the online supplement. (B) The values of μ_{EI}^{-1} and τ (in units of days) for these 300 parameter estimates. Points with high log likelihood, within 5 log units of the maximum, are shown as shaded circles (green in the online version). The remaining points are represented by black crosses. The online version of this figure is in color.

that the total latent period is more precisely identified than the separate components. The top left cluster in Figure 3B has mosquito latent period close to biologically plausible values of around two weeks; the bottom right cluster has human latent period close to biologically plausible values of around one week. The highlighted points in Figure 3B reveal that the top left cluster is associated with the higher likelihood values. For both clusters, the total latent period is larger than one might expect from a biological perspective. Recently, it has been proposed that upward bias in estimated latent and infectious periods can be a consequence of employing models without spatial structure for large populations (He, Ionides, and King 2010). In addition, the system delays due to mosquito population dynamics and the response to rainfall have simplified representations in our model that could potentially bias the estimation of other latent periods. What our analysis can legitimately investigate is the effective value of these parameters consistent with the use of population-level aggregate models to describe aggregate data. Public health is concerned with population-level quantities, and so the relationship between population and individual level models is a relevant topic for investigation in its own right.

A simple, but valuable, diagnostic for the specification of a mechanistic model is to compare the goodness of fit with standard nonmechanistic statistical models. One can argue that part of the point of fitting a mechanistic model to data is to discover which aspects of the data are not captured by a model describing current scientific knowledge about the system under investigation. Somewhat equivalently, one might understand that requiring a model to have scientific interpretability may lead to a cost in terms of the ability to match data statistically. In this sense, it may not be a scientific goal to achieve a level of fit comparable to flexible statistical models which do not seek scientific interpretability. On the other hand, to carry out formal hypothesis tests, or to interpret parameter estimates and their uncertainty, it is helpful if the model can be shown to give an adequate statistical fit to the data. In Table 1, we include as a benchmark comparison a model in which $\{\log(y_n + 1), n = 1, \dots, N\}$ is supposed to follow a Gaussian SARIMA specification. The large number of additional parameters in the mechanistic models appears to be justified relative to this log-SARIMA model, according to the AIC criterion. Log-SARIMA models are theoretically appealing as simple models for disease transmission, since (in common with many other biological populations) the *Plasmodium* demonstrates annual cycles of abundance which consist approximately of a period of exponential growth followed by

a period of exponential decay. As another benchmark, we included the rainfall covariate Z_t into the log-SARIMA model (via the ARMAX framework; Shumway and Stoffer 2006), also reported in Table 1. The improvement in model fit from including the covariate is comparable, in terms of units of log likelihood, to the improvement seen in the VSEIR and VS²EI² models.

From Table 1, we see that all the four mechanistic models analyzed beat the benchmark nonmechanistic log-SARIMA model by a large margin of AIC. Having established that these models are adequate statistical explanations of the data, we compare these models amongst each other. Likelihoods for both the VS²EI² model and the simpler VSEIR submodel (described in the caption to Figure 1) improve significantly when the rainfall covariate is used ($p < 0.001$ for the likelihood ratio test, using a chi-square approximation on one degree of freedom). After concluding that inclusion of rainfall does indeed help to describe malaria dynamics, we proceed to compare the VSEIR and the VS²EI² models, both including rainfall. These two models have different numbers of parameters and we can compare their Akaike Information Criterion (AIC) values, which favors the VS²EI² model with rainfall. Since these two models are nested, one can also carry out a likelihood ratio test of the null hypothesis that the data follow the VSEIR model ($p < 0.001$, chi-square test on 5 degrees of freedom). The nesting is nonstandard (e.g., when $\mu_{I_2S_2} \rightarrow \infty$ the initial value $[I_2]_0$ becomes undefined), however, the chi-square test is expected to be conservative in such situations (Anisimova, Bielawski, and Yang 2001). This comparison is evidence for the value of incorporating characteristic aspects of the human immune response to malaria into models used for time series analysis. However, models based on simpler SEIR descriptions of human immunity will continue to be central to the study of disease dynamics, and our results also support a position that the VSEIR model is not entirely discredited. It produces parameter estimates which are qualitatively similar to the VS²EI² model, and its log likelihood is much closer to that of the VS²EI² than to the log-SARIMA benchmark. To understand the relative strengths and weaknesses of different models, one pertinent question to consider is which parts of the time series are better explained by each model. In Figure 4 we plot the difference of the conditional log likelihoods of the VS²EI² model with rainfall and the VSEIR model with rainfall, at each point in time. We note that during many of the epidemics, most notably in the fall of 1989, 1990, 1992, 1994, and 1997, the simpler VSEIR model fits the data better as the epidemic approaches its peak. Predicting the

Table 1. A likelihood-based comparison of the fitted models

Model	Log likelihood (ℓ)	p	AIC
VSEIR without rainfall	-1275.0	19	2588.0
VSEIR with rainfall	-1265.0	20	2570.0
VS ² EI ² without rainfall	-1261.1	24	2570.2
VS ² EI ² with rainfall	-1251.0	25	2552.0
Log-SARIMA (1, 0, 1) \times (1, 0, 1) ₁₂ without rainfall	-1329.0	6	2670.0
Log-SARIMA (1, 0, 1) \times (1, 0, 1) ₁₂ with rainfall	-1322.6	7	2659.2

NOTE: Corresponding point estimates are presented in Table 3. The column labeled p corresponds to the number of estimated parameters, including unknown initial conditions. Parameters which were not estimated are documented in Table 2. AIC is computed as $AIC = -2\ell + 2p$.

Table 2. List of symbols used in the article with a description and units

Symbol	Brief description	Unit	Fixed value
μ_{XY}	Per-capita transition rate from X to Y ; $X, Y \in \{S_1, S_2, E, I_1, I_2\}$	yr^{-1}	—
$[X]_0$	Initial fraction in compartment X ; $X \in \{S_1, S_2, E, I_1, I_2\}$	—	—
$[\lambda_i]_0$	Initial values for the latent force of infection ($i = 1, \dots, k$)	—	—
τ	Mean development delay for mosquitoes	yr	—
σ	Standard deviation of the process noise	$\text{yr}^{1/2}$	—
ρ	Reporting fraction	—	—
q	Relative infectivity of partially immune individuals	—	—
c	Coefficient of reinfection with clinical immunity	—	—
k	Shape parameter for the delay development kernel for mosquitoes	—	2
ψ	Dispersion parameter of the observation noise	—	—
n_s	Number of splines describing seasonality	—	6
β_i	Spline coefficients, for $i = 1, \dots, n_s$	—	—
$\bar{\beta}$	Dimensionality constant	yr^{-1}	1
β	Coefficient of climate (rainfall) covariate	—	—
u	Window for rainfall to affect transmission	mo	5
v	Threshold for integrated rainfall	mm	200
$1/\delta$	Average life expectancy	yr	50
Δ	Time step for stochastic Euler integration	day	1

NOTE: Some parameters were not estimated for the analysis presented in this article, and the last column gives their fixed values. A justification of the choice $k = 2$ is provided in the online supplement. The values of estimated parameters are give in Table 3.

peak of an epidemic is of particular public health interest, as it determines the maximum case burden experienced by the health care system. The more complex VS^2EI^2 model, which fits the data better overall, may have little or no advantage for this specific purpose.

When building mechanistic dynamic models for biological systems, there is a temptation to include as much biological detail as the available data will support. A price for this is that certain combinations of parameters may be weakly identified by the data. However, we can focus on conclusions which are robust to identifiability issues. For example, the model comparison via log likelihoods in Table 1 is valid despite any potential lack of identifiability. Although point estimates provide a convenient summary of a fitted model, they should be interpreted with caution in the presence of weak identifiability.

Confidence intervals, which become wide as the statistical evidence becomes weak, are more appropriate for drawing scientific conclusions. The profile likelihood for the reporting rate in Figure 5 shows that, without making any specific assumptions on the the values of the 25 parameters estimated, there is evidence that the effective reporting rate is less than 2.5%. There is general agreement that malaria is substantially under-reported in Southeast Asia (Snow et al. 2005) and a study in the city of Ahmedabad, Gujarat, found a reporting rate of 10% (Yadav et al. 2003). Much of the remaining discrepancy could be explained by a recent suggestion, based on a sensitive polymerase chain reaction diagnostic analysis in an epidemic malaria region of the East African highlands, that microscopy techniques may fail to detect two thirds of asymptomatic *Plasmodium* infections (Baliraine et al. 2009). There is potential for asymptomatic in-

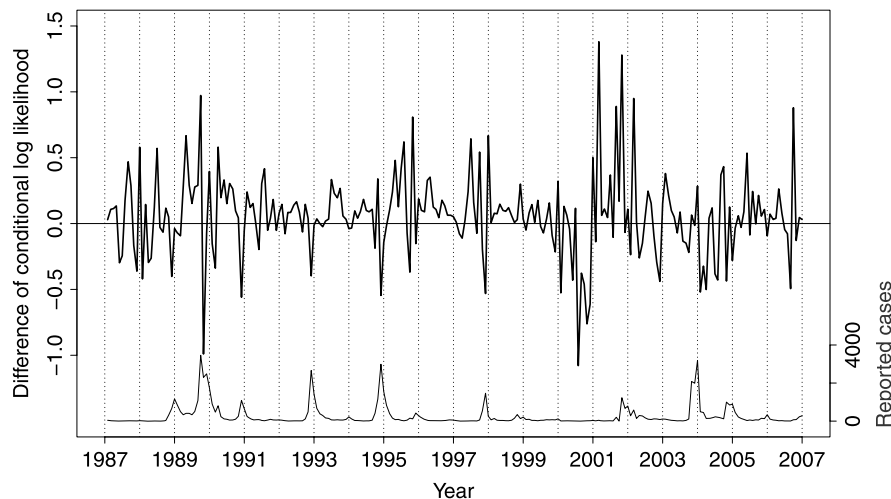


Figure 4. Difference between the conditional log likelihood of y_n given y_1, \dots, y_{n-1} for the VS^2EI^2 model with rainfall and the VSEIR model with rainfall, plotted against time (bold line). For comparison, reported malaria cases in Kutch are also shown (thin line).

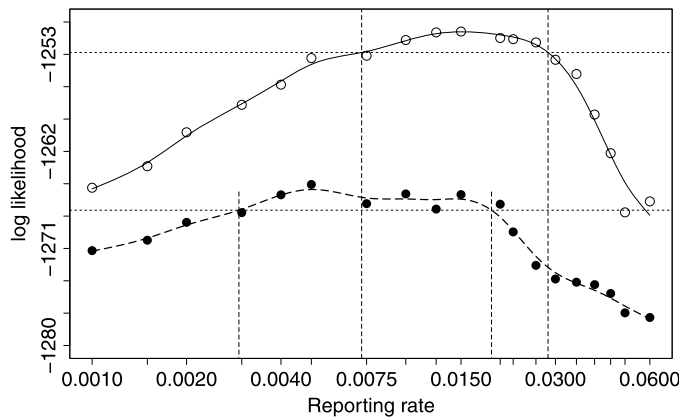


Figure 5. Profile likelihood plot for the reporting rate (ρ) for the VS^2EI^2 model with rainfall (solid line) and the VSEIR model with rainfall (broken line). The profile is estimated via fitting a smooth curve through Monte Carlo evaluations shown as open circles (VS^2EI^2) and filled circles (VSEIR). The dashed vertical lines construct 95% confidence intervals using a cutoff of 1.92 (from a chi-square approximation; e.g., [Barndorff-Nielsen and Cox 1994](#)) for this estimated profile.

fections to play important dynamic roles, which can be hard to identify ([King et al. 2008](#)); for example, there could be an epidemiological role for boosted immunity due to mild infections that occur at blood parasite levels too low to be detected by standard field investigations. One cannot at this point rule out the possibility that the low estimated reporting rate could be an artifact due to unmodeled population inhomogeneity, or some other shortcoming of the model. Resolving such questions is beyond the scope of this article. The statistical interpretation, however, is more clearcut: Any attempt to learn about malaria via fitting epidemiological models of the type constructed here must take into account the discovery that unconventionally low reporting rates may give superior explanation of the data.

The chi-square cutoffs used to generate the approximate confidence intervals in Figure 5, and the likelihood ratio tests used to interpret Table 1, have an asymptotic justification ([Barndorff-Nielsen and Cox 1994](#)). We investigated the actual coverage probability of the interval in Figure 5 at the maximum likelihood estimate (MLE), for the VS^2EI^2 model, by simulating from the model at the MLE parameter values and reconstructing a profile on reporting rate for each simulation. For 200 such simulated profiles, we found an empirical coverage probability of 96.5% which suggests that the actual coverage probability is close to its nominal value. The discrepancy between actual and nominal coverage may be larger for other parameters (see Figure 6 for an example). These simulations could in principle be used to determine a profile likelihood cutoff that has exactly the desired nominal coverage at the MLE; we do not investigate this further here since the likelihood ratios used to support our substantial conclusions are sufficiently large compared to the asymptotically justified significance cutoffs to be statistically unambiguous.

Many parameter estimates have large statistical uncertainty (Table 3, last two columns). One could investigate whether fixing some parameters at previously published scientific values helps to identify some other parameters. Conclusions from such

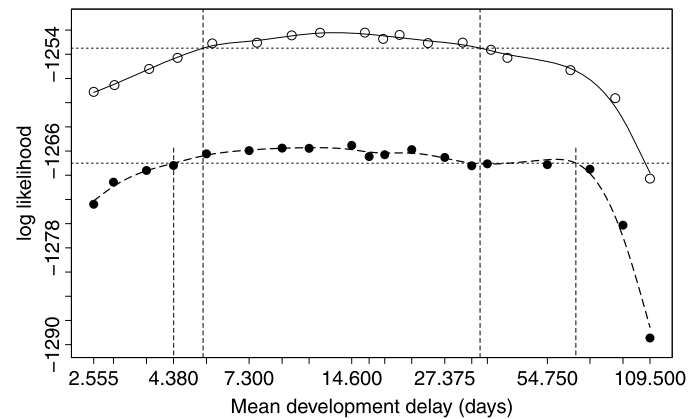


Figure 6. Profile likelihood plot for the mean development delay time of mosquitoes (τ) for the VS^2EI^2 model with rainfall (solid line) and the VSEIR model with rainfall (broken line). The dashed vertical lines construct approximate 95% confidence intervals, as described in Figure 5. The empirical coverage probability for 200 simulations from the MLE of the VS^2EI^2 model was 82.5%, indicating that the actual uncertainty is somewhat larger than these intervals indicate.

an analysis should be made cautiously, since the variability and complexity of biological systems means that it is typically difficult to know to what extent previous investigations are indeed quantitatively relevant for the current model and data. This consideration would similarly complicate the development of a scientifically informed prior distribution, if one were to investigate a Bayesian approach.

Adding additional parameters to a model does not necessarily result in more weakly identified parameter estimates, particularly when the extended model provides substantial improvement in fit. Figure 6 provides one such example, where the *Plasmodium* development delay τ is more precisely estimable in the larger VS^2EI^2 model. Further, in the VS^2EI^2 model the parameter values which are consistent with the data are closer to the biological interpretation as a development delay—directly measured mean development times are around two weeks in this context. It could be nothing but a happy accident that, in this case, a model estimated on population data happens to match an individual-level biological interpretation. Given all the simplifications necessarily involved in the modeling process, it is hard to be sure that this parameter describes the biological interpretation in the strong sense that manipulation of the development time would affect the system only through the estimated value of τ . Since development time is a well-studied function of temperature, in principle one could investigate this by seeing whether building this dependence into the model improves its explanation of the data. However, even without insisting on such a strong interpretation, when the data and the model and the desired biological interpretation are all mutually consistent then the model becomes validated as a conceptual tool for understanding the biological system.

4. DISCUSSION

One of the clearest scientific conclusions from our data analysis is that rainfall variability does indeed have a detectable effect on malaria dynamics in Kutch, even once one controls for seasonality and nonlinear dynamic effects of the force of infection and immunity. This is a contribution to the debate on

Table 3. Estimated model parameters

	VSEIR without rain	VS ² EI ² without rain	VSEIR with rain	VS ² EI ² with rain	Confidence interval
$\mu_{I_1S_2}$	13.587	—	39.021	—	(—, —)
$\mu_{S_2S_1}$	0.116	0.230	5.657	0.334	(0.067, 3.270)
μ_{EI_1}	7.301	7.408	10.480	8.902	(8.885, 17.277)
$\mu_{I_1I_2}$	—	11.544	—	5.511	(3.218, ∞)
$\mu_{I_2S_2}$	—	0.004	—	0.035	(0, 0.073)
$\mu_{I_1S_1}$	—	2.320	—	6.563	(0, ∞)
β_1	−0.076	−2.469	1.242	1.201	(−4.819, 4.109)
β_2	1.287	2.001	3.590	2.088	(−0.153, 6.616)
β_3	4.446	4.227	3.906	3.866	(1.874, 6.939)
β_4	2.868	2.786	3.747	2.808	(1.092, 6.042)
β_5	6.709	6.534	5.742	5.996	(4.695, 9.749)
β_6	6.319	7.080	4.803	5.333	(3.912, 8.287)
τ	0.025	0.022	0.033	0.030	(0.015, 0.084)
σ	0.347	0.309	0.225	0.243	(0.162, 0.259)
ρ	0.022	0.030	0.005	0.015	(0.007, 0.025)
$q \times 10^4$	—	4.763	—	9.424	(0.100, 48.102)
ψ	0.384	0.390	0.390	0.395	(0.365, 0.445)
β	—	—	0.489	0.512	(0.270, 0.765)
$[S_1]_0$	0.494	0.164	0.956	0.138	(0.001, 0.900)
$[S_2]_0$	0.505	0.765	0.038	0.775	(0.276, 0.900)
$[E]_0$	0.003	0.002	0.014	0.004	(0.003, 0.009)
$[I_1]_0$	0.011	0.002	0.002	0.002	(0, 0.087)
$[I_2]_0$	—	0.067	—	0.080	(0, 0.754)
$[\lambda_1]_0 \times 10$	0.079	0.133	0.189	0.171	(0, ∞)
$[\lambda_2]_0 \times 10$	0.050	0.045	0.058	0.061	(0, ∞)
c	—	0.004	—	0.010	(0.001, 0.067)

NOTE: The columns marked ‘without rain’ correspond to maximum likelihood point estimates under the constraint $\beta = 0$. The last two columns give the lower and upper bounds for approximate 95% confidence intervals for the VS²EI² model with rainfall, derived from profile likelihood computations as shown in Figures 5 and 6; values of 0 and ∞ correspond to confidence intervals extending to the boundary of the parameter space. These models were also analyzed by Laneri et al. (2010).

the role of climate variability in malaria transmission, which has previously been lacking such an analysis (Hay et al. 2002; Pascual et al. 2006, 2008; Briët et al. 2008). We have studied just one district here, in order to focus on the statistical principles behind our analysis. The statistical approach presented will facilitate similar investigations of other regions with endemic and epidemic malaria. Given geographical differences in mosquito species, social and agricultural practices, and many relevant ecological variables, one should however be cautious about extrapolating our quantitative results.

At the African summit on Roll Back Malaria in April 2000, 44 leaders of affected countries signed the Abuja declaration. One of the requirements of this declaration was that malaria epidemics should be detected, and effective control measures implemented, within two weeks. In practice, this timeline necessitates the use of malaria forecasts. Two major components of such a forecast should be measures of environmental suitability for transmission and the extent of residual immunity from previous epidemics. Seasonal rainfall forecasts (Bouma and van der Kaay 1994, 1996) and satellite observations (Thomson et al. 2006) may have a role to play. Indeed, local rainfall was used as the primary component of epidemic malaria forecasts published for the semi-arid Punjab in the early part of the 20th century (Swaroop 1949). Our results confirm that local rainfall acts at a sufficient lag to be a simple and useful predictor. The

utility of our models for forecasting was investigated by Laneri et al. (2010).

This case study has demonstrated a statistical framework for likelihood-based inference using nonlinear, partially observed, multivariate Markov process models. Reasons for considering frameworks other than likelihood-based inference include (i) a preference for Bayesian analysis; (ii) a concern that likelihood-based methodology does not have theoretical optimality guarantees for finite-sample inference; (iii) the possibility that alternative methodology might be more computationally convenient, whether or not it has comparable statistical efficiency. On the other hand, likelihood-based inference is a widely employed paradigm that has been applied successfully in many situations. This provides motivation for extending its applicability to classes of complex dynamic models that are playing increasing roles in ecology, epidemiology and elsewhere. Much recent work in the area of inference for POMP models has followed the Bayesian paradigm (e.g., Boys, Wilkinson, and Kirkwood 2008; Cauchemez and Ferguson 2008; McKinley, Cook, and Deardon 2009; Toni et al. 2009; Andrieu, Doucet, and Holenstein 2010; Wilkinson 2011). We have demonstrated that maximum likelihood methodology can be a computationally viable alternative to these Bayesian approaches, in addition to being readily applicable due to the plug-and-play property. The analysis presented in this article is consistent with a study by Liu et al. (2009) which reported computational advantages for

adopting a maximum likelihood approach over Bayesian methods for inference on complex phylogenetic models. Regardless of one's opinion on the epistemological value of asserting a prior distribution on unknown parameters, there may sometimes be computational advantages to exploring the likelihood surface rather than a posterior distribution.

Other vector-borne diseases, such as dengue and leishmaniasis, lead to statistical considerations and challenges similar to those for malaria. In a wider context, disease systems exemplify the issues at stake in developing an understanding of ecological processes from available time-series data (Bjornstad and Grenfell 2001). Quantitative understanding of ecosystems has growing importance as mankind is increasingly responsible for managing the biological resources of the planet. The broad scope of these responsibilities will continue to drive further developments in statistical methodology and data analysis.

SUPPLEMENTARY MATERIALS

Supplementary text: Pdf file presenting the implementation of iterated filtering used for the calculations in this article and two additional figures. (supplement.pdf)

[Received May 2010. Revised February 2011.]

REFERENCES

- Anderson, R. M., and May, R. M. (1991), *Infectious Diseases of Humans*, Oxford: Oxford University Press. [443]
- Andrieu, C., Doucet, A., and Holenstein, R. (2010), "Particle Markov Chain Monte Carlo," *Journal of the Royal Statistical Society, Ser. B*, 72, 269–342. [441,449]
- Anisimova, M., Bielawski, J. P., and Yang, Z. (2001), "Accuracy and Power of the Likelihood Ratio Test in Detecting Adaptive Molecular Evolution," *Molecular Biology and Evolution*, 18, 1585–1592. [446]
- Applebaum, D. (2004), "Lévy Processes: From Probability to Finance and Quantum Groups," *Notices of the American Mathematical Society*, 51, 1336–1347. [441,443]
- Aron, J. L., and May, R. M. (1982), "The Population Dynamics of Malaria," in *The Population Dynamics of Infectious Diseases*, ed. R. M. Anderson, London, U.K.: Chapman & Hall, pp. 139–179. [441,443]
- Artavanis-Tsakonas, K., Tongren, J. E., and Riley, E. M. (2003), "The War Between the Malaria Parasite and the Immune System: Immunity, Immunoregulation and Immunopathology," *Clinical and Experimental Immunology*, 133, 145–152. [442]
- Baliraine, F., Afrane, Y., Ameny, D., Bonizzoni, M., Menge, D., Zhou, G., Zhong, D., Vardo-Zalik, A., Githeko, A., and Yan, G. (2009), "High Prevalence of Asymptomatic *Plasmodium falciparum* Infections in a Highland Area of Western Kenya: A Cohort Study," *The Journal of Infectious Diseases*, 200, 66–74. [447]
- Barndorff-Nielsen, O. E., and Cox, D. R. (1994), *Inference and Asymptotics*, London: Chapman & Hall. [448]
- Bartlett, M. S. (1960), *Stochastic Population Models in Ecology and Epidemiology*, New York: Wiley. [441]
- Bhadra, A. (2010), Discussion of "Particle Markov Chain Monte Carlo Methods," by C. Andrieu, A. Doucet, and R. Holenstein, *Journal of the Royal Statistical Society, Ser. B*, 72, 314–315. [441]
- Bjornstad, O. N., and Grenfell, B. T. (2001), "Noisy Clockwork: Time Series Analysis of Population Fluctuations in Animals," *Science*, 293, 638–643. [444,450]
- Bouma, M. J., and van der Kaay, H. J. (1994), "Epidemic Malaria in India and the El Niño Southern Oscillation," *The Lancet*, 344, 1638–1639. [444,449]
- (1996), "The El Niño Southern Oscillation and the Historic Malaria Epidemics on the Indian Subcontinent and Sri Lanka: An Early Warning System for Future Epidemics?" *Tropical Medicine and International Health*, 1, 86–96. [449]
- Boys, R. J., Wilkinson, D. J., and Kirkwood, T. B. L. (2008), "Bayesian Inference for a Discretely Observed Stochastic Kinetic Model," *Statistica Sinica*, 18, 125–135. [449]
- Bretó, C., He, D., Ionides, E. L., and King, A. A. (2009), "Time Series Analysis via Mechanistic Models," *The Annals of Applied Statistics*, 3, 319–348. [441,443,445]
- Brët, O., Vounatsou, P., Gunawardena, D., Galappaththy, G., and Amerasinghe, P. (2008), "Temporal Correlation Between Malaria and Rainfall in Sri Lanka," *Malaria Journal*, 7, 77. [449]
- Cappé, O., Moulines, E., and Rydén, T. (2005), *Inference in Hidden Markov Models*, New York: Springer. [441]
- Cauchemez, S., and Ferguson, N. M. (2008), "Likelihood-Based Estimation of Continuous-Time Epidemic Models From Time-Series Data: Application to Measles Transmission in London," *Journal of the Royal Society Interface*, 5, 885–897. [441,449]
- Cauchemez, S., Valleron, A., Boëlle, P., Flahault, A., and Ferguson, N. M. (2008), "Estimating the Impact of School Closure on Influenza Transmission From Sentinel Data," *Nature*, 452, 750–754. [441]
- Chitnis, N., Cushing, J. M., and Hyman, J. M. (2006), "Bifurcation Analysis of a Mathematical Model for Malaria Transmission," *SIAM Journal on Applied Mathematics*, 67, 24–45. [443]
- Dietz, K., Molineaux, L., and Thomas, A. (1974), "A Malaria Model Tested in the African Savannah," *Bulletin of the World Health Organization*, 50, 347–357. [441]
- Earn, D. J. D., Rohani, P., Bolker, B. M., and Grenfell, B. T. (2000), "A Simple Model for Complex Dynamical Transitions in Epidemics," *Science*, 287, 667–670. [444]
- Ellner, S. P., Bailey, B. A., Bobashev, G. V., Gallant, A. R., Grenfell, B. T., and Nychka, D. W. (1998), "Noise and Nonlinearity in Measles Epidemics: Combining Mechanistic and Statistical Approaches to Population Modeling," *American Naturalist*, 151, 425–440. [441]
- Filipe, J. A. N., Riley, E. M., Drakeley, C. J., Sutherland, C. J., and Ghani, A. C. (2007), "Determination of the Processes Driving the Acquisition of Immunity to Malaria Using a Mathematical Transmission Model," *PLoS Computational Biology*, 3, e255. [441,442]
- Finkenstädt, B. F., and Grenfell, B. T. (2000), "Time Series Modelling of Childhood Diseases: A Dynamical Systems Approach," *Applied Statistics*, 49, 187–205. [441]
- Greenwood, B. (2009), "Can Malaria Be Eliminated?" *Transactions of the Royal Society of Tropical Medicine and Hygiene*, 103, S2–S5. [440]
- Gupta, S., Trenholme, K., Anderson, R., and Day, K. (1994), "Antigenic Diversity and the Transmission Dynamics of *Plasmodium falciparum*," *Science*, 263, 961–963. [441,443]
- Hay, S. I., Cox, J., Rogers, D. J., Randolph, S. E., Stern, D. I., Shanks, G. D., Myers, M. F., and Snow, R. W. (2002), "Climate Change and the Resurgence of Malaria in the East African Highlands," *Nature*, 415, 906–909. [442,449]
- Hay, S. I., Guerra, C. A., Tatem, A. J., Atkinson, P. M., and Snow, R. W. (2005), "Tropical Infectious Diseases: Urbanization, Malaria Transmission and Disease Burden in Africa," *Nature Reviews Microbiology*, 3, 81–90. [440]
- He, D., Ionides, E. L., and 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. [441,445,446]
- Ionides, E. L., Bretó, C., and King, A. A. (2006), "Inference for Nonlinear Dynamical Systems," *Proceedings of the National Academy of Sciences of the USA*, 103, 18438–18443. [441,445]
- Jacod, J. (2004), "The Euler Scheme for Lévy Driven Stochastic Differential Equations: Limit Theorems," *The Annals of Probability*, 32, 1830–1872. [441,443]
- Keeling, M., and Rohani, P. (2008), *Modeling Infectious Diseases in Humans and Animals*, Princeton, NJ: Princeton University Press. [443]
- Kendall, B. E., Ellner, S. P., McCauley, E., Wood, S. N., Briggs, C. J., Murdoch, W. W., and Turchin, P. (2005), "Population Cycles in the Pine Looper Moth: Dynamical Tests of Mechanistic Hypotheses," *Ecological Monographs*, 75, 259–276. [441]
- King, A. A., Ionides, E. L., Bretó, C., Ellner, S., and Kendall, B. (2009), "pomp: Statistical Inference for Partially Observed Markov Processes," R package, available at www.r-project.org. [445]
- King, A. A., Ionides, E. L., Pascual, M., and Bouma, M. J. (2008), "Inapparent Infections and Cholera Dynamics," *Nature*, 454, 877–880. [445,448]
- Kiszewski, A. E., and Teklehaimanot, A. (2004), "A Review of the Clinical and Epidemiologic Burdens of Epidemic Malaria," *American Journal of Tropical Medicine and Hygiene*, 71, 128–135. [441,444]
- Klein, E., Smith, D., Boni, M., and Laxminarayan, R. (2008), "Clinically Immune Hosts as a Refuge for Drug-Sensitive Malaria Parasites," *Malaria Journal*, 7, 67. [441–443]
- Koella, J., and Antia, R. (2003), "Epidemiological Models for the Spread of Anti-Malarial Resistance," *Malaria Journal*, 2, 3. [441,443]
- Kumar, A., Valecha, N., Jain, T., and Dash, A. P. (2007), "Burden of Malaria in India: Retrospective and Prospective View," *American Journal of Tropical Medicine and Hygiene*, 77, 69–78. [440,444]
- Laner, K., Bhadra, A., Ionides, E. L., Bouma, M., Yadav, R., Dhiman, R., and Pascual, M. (2010), "Forcing versus Feedback: Epidemic Malaria and Monsoon Rains in NW India," *PLoS Computational Biology*, 6, e1000898. [442,445,449]

- Liu, J., and West, M. (2001), "Combining Parameter and State Estimation in Simulation-Based Filtering," in *Sequential Monte Carlo Methods in Practice*, eds. A. Doucet, N. de Freitas, and N. J. Gordon, New York: Springer, pp. 197–224. [441]
- Liu, K., Raghavan, S., Nelesen, S., Linder, C. R., and Warnow, T. (2009), "Rapid and Accurate Large-Scale Coestimation of Sequence Alignments and Phylogenetic Trees," *Science*, 324, 1561–1564. [449]
- Lloyd, A. L. (2001), "Realistic Distributions of Infectious Periods in Epidemic Models: Changing Patterns of Persistence and Dynamics," *Theoretical Population Biology*, 60, 59–71. [443]
- Macdonald, G. (1957), *The Epidemiology and Control of Malaria*, London: Oxford University Press. [440,443]
- McFadden, D. (1989), "A Method of Simulated Moments for Estimation of Discrete Response Models Without Numerical Integration," *Econometrica*, 57, 995–1026. [441]
- McKenzie, F. E., and Bossert, W. H. (2005), "An Integrated Model of *Plasmodium falciparum* Dynamics," *Journal of Theoretical Biology*, 232, 411–426. [441,443]
- McKenzie, F. E., and Samba, E. M. (2004), "The Role of Mathematical Modeling in Evidence-Based Malaria Control," *American Journal of Tropical Medicine and Hygiene*, 71, 94–96. [440]
- McKinley, T., Cook, A. R., and Deardon, R. (2009), "Inference in Epidemic Models Without Likelihoods," *The International Journal of Biostatistics*, 5, article 24. [441,449]
- Packard, R. M. (2007), *The Making of a Tropical Disease: A Short History of Malaria*, Baltimore, MD: Johns Hopkins University Press. [440,441]
- Pascual, M., Ahumada, J., Chaves, L. F., Rodo, X., and Bouma, M. (2006), "Malaria Resurgence in East African Highlands: Temperature Trends Revisited," *Proceedings of the National Academy of Sciences of the USA*, 103, 5829–5835. [442,449]
- Pascual, M., Cazelles, B., Bouma, M., Chaves, L., and Koelle, K. (2008), "Shifting Patterns: Malaria Dynamics and Rainfall Variability in an African Highland," *Proceedings of the Royal Society B: Biological Sciences*, 275, 123–132. [444,449]
- Protter, P., and Talay, D. (1997), "The Euler Scheme for Lévy Driven Stochastic Differential Equations," *The Annals of Probability*, 25, 393–423. [441,443]
- Roberts, G. O., and Stramer, O. (2001), "On Inference for Partially Observed Nonlinear Diffusion Models Using the Metropolis–Hastings Algorithm," *Biometrika*, 88, 603–621. [441]
- Roberts, L., and Enserink, M. (2007), "Malaria: Did They Really Say... Eradication?" *Science*, 318, 1544–1545. [440]
- Ross, R. (1911), *The Prevention of Malaria*, London: John Murray. [440]
- Sharma, V. P. (1996), "Re-Emergence of Malaria in India," *Indian Journal of Medical Research*, 103, 26–45. [440]
- Shumway, R. H., and Stoffer, D. S. (2006), *Time Series Analysis and Its Applications* (2nd ed.), New York: Springer. [446]
- Sisson, S. A., Fan, Y., and Tanaka, M. M. (2007), "Sequential Monte Carlo Without Likelihoods," *Proceedings of the National Academy of Sciences of the USA*, 104, 1760–1765. [441]
- Snow, R. W., Guerra, C. A., Noor, A. M., Myint, H. Y., and Hay, S. I. (2005), "The Global Distribution of Clinical Episodes of *Plasmodium falciparum* Malaria," *Nature*, 434, 214–217. [440,447]
- Swaroop, S. (1949), "Forecasting of Epidemic Malaria in the Punjab, India," *American Journal of Tropical Medicine and Hygiene*, 51-29, 1–17. [444, 449]
- Thomson, M. C., Doblas-Reyes, F. J., Mason, S. J., Hagedorn, S. J., Phindela, T., Morse, A. P., and Palmer, T. N. (2006), "Malaria Early Warnings Based on Seasonal Climate Forecasts From Multi-Model Ensembles," *Nature*, 439, 576–579. [449]
- Toni, T., Welch, D., Strelkowa, N., Ipsen, A., and Stumpf, M. P. (2009), "Approximate Bayesian Computation Scheme for Parameter Inference and Model Selection in Dynamical Systems," *Journal of the Royal Society Interface*, 6, 187–202. [449]
- Warrell, D. A., and Gilles, H. A. (eds.) (2002), *Essential Malariology*, London: Arnold. [442]
- Wilkinson, D. J. (2011), "Parameter Inference for Stochastic Kinetic Models of Bacterial Gene Regulation: A Bayesian Approach to Systems Biology" (with discussion), in *Bayesian Statistics 9*, eds. J. M. Bernardo, M. J. Bayarri, J. O. Berger, A. P. Dawid, D. Heckerman, A. F. M. Smith, and M. West, Oxford: Oxford University Press, to appear. [441,449]
- Wood, S. N. (2010), "Statistical Inference for Noisy Nonlinear Ecological Dynamic Systems," *Nature*, 466, 1102–1104. [441]
- World Health Organization (2008), *World Malaria Report*, Geneva: WHO Press. [440]
- Yadav, R. S., Bhatt, R. M., Kohli, V. K., and Sharma, V. P. (2003), "The Burden of Malaria in Ahmedabad City, India: A Retrospective Analysis of Reported Cases and Deaths," *Annals of Tropical Medicine and Parasitology*, 97, 793–802. [447]