

Innovations in Likelihood-Based Inference for State Space Models

by

Jesse Wheeler

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Statistics)
in The University of Michigan
2025

Doctoral Committee:

Professor Edward L. Ionides, Chair
Professor Aaron A. King
Assistant Professor Jeffrey Regier
Professor Kerby Shedden

Jesse Wheeler

jeswheel@umich.edu

ORCID iD: [0000-0003-3941-3884](https://orcid.org/0000-0003-3941-3884)

© Jesse Wheeler 2025

DEDICATION

Dedicated to my family: past, present, and future.

ACKNOWLEDGEMENTS

I am deeply grateful to all of the mentors in my life who have helped me get to this point. I am particularly grateful for my advisor, Dr. Edward L. Ionides, who has provided invaluable support for me throughout my academic journey and has helped shape my beliefs about statistics and higher education. I wish to also thank the other members of my committee who have provided their advise and mentorship. Dr. Aaron A. King

Most importantly I would like to thank my spouse, Haylee Wheeler, who has been incredibly supportive and encouranging though my entire academic career. The progress that I have made as a scholar would not have been possible without her. You have always been willing to listen to my complaints and struggles, and been with me to celebrate life's victories.

TABLE OF CONTENTS

DEDICATION	ii
ACKNOWLEDGEMENTS	iii
LIST OF FIGURES	vi
LIST OF TABLES	vii
LIST OF APPENDICES	viii
ABSTRACT	ix

CHAPTER

1 Introduction	1
1.1 Defining State Space Models and Notation	2
1.2 Overview of remaining chapters	3
2 Likelihood Based Inference for ARMA Models	5
2.1 Introduction	5
2.2 Maximum Likelihood for ARMA Models	7
2.2.1 A Novel Multi-start Algorithm	9
2.2.2 Simulation Studies	12
2.3 Annual Depths of Lake Michigan	18
2.3.1 Parameter uncertainty	20
2.4 Discussion	22
3 Informing policy via dynamic models: Cholera in Haiti	24
3.1 Introduction	24
3.2 Materials and methods	27
3.2.1 Mechanistic models for disease modeling	27
3.2.2 Model Fitting	39
3.2.3 Model Diagnostics	47
3.2.4 Corroborating Fitted Models with Scientific Knowledge	50
3.3 Results	52
3.3.1 Forecasts	52
3.4 Discussion	55

3.4.1 Reproducibility and Extendability	57
4 Conclusion	59
APPENDICES	60
BIBLIOGRAPHY	90

LIST OF FIGURES

FIGURE

2.1	Multi-node MA(1) models	9
2.2	Proportion of improved models	13
2.3	Confidence interval coverage	15
2.4	AIC table improvements	17
2.5	Average depth of Lake Michigan-Huron from 1860-2014.	18
2.6	Evidence for an AR(1) model for the Lake Michigan-Huron data.	21
2.7	ARMA polynomial roots	21
3.1	Weekly cholera cases	26
3.2	Trend in transmission rate confidence interval.	42
3.3	Simulations from Model 1 compared to reported cholera cases.	43
3.4	Simulated trajectory of Model 2	45
3.5	Simulations from Model 3 compared to reported cholera cases.	47
3.6	Seasonality of Model 1 transmission rates and weekly rainfall.	51
3.7	Probability of cholera elimination under Models 1 and 3.	54
A.1	Uniform resampling of model parameters.	61
B.1	Flow diagram of Haiti-cholera model (Model 1)	63
B.2	Flow diagram of Haiti-cholera model (Model 2)	64
B.3	Flow diagram of Haiti-cholera model (Model 3)	65
B.4	MCAP confidence intervals for Model 1 parameters.	69
B.5	MCAP confidence intervals for Model 2 parameters.	71
B.6	MCAP confidence intervals for Model 3 parameters.	72
B.7	Replicating Model 1 epidemic phase parameter distributions	74
B.8	Replicating Model 1 endemic phase parameter distributions	75
B.9	Simulations from replicated Model 1	76
B.10	Model 2 trajectories.	77
B.11	Simulations from replicated Model 3	79
B.12	Model 3 unit log-likelihoods.	81
B.13	Conditional log-likelihoods without hurricane adjustment	81
B.14	Conditional log-likelihoods with hurricane adjustment	82
B.15	Log-likelihoods of Model 3 for each department compared to the corresponding benchmark model after adding and estimating parameters related to Hurricane Matthew.	83
B.16	Estimated susceptible population over time.	85

LIST OF TABLES

TABLE

2.1	AIC values of ARMA model for Lake Michigan-Huron	19
2.2	Parameter values of an ARMA model fit to Lake Michigan-Huron data.	20
3.1	Parameter estimates of the cholera models.	29
3.2	Parameter translations to original model definitions	30
3.3	AIC values for cholera models and benchmarks	40
B.1	Model 1 parameter estimates and confidence intervals.	70
B.2	Model 2 parameter estimates and confidence intervals.	70
B.3	Model 3 parameter estimates and confidence intervals.	71

LIST OF APPENDICES

A Appendix for Chapter 2	60
B Appendix for Chapter 3	62
C Example Appendix 01	87
D Example Appendix 02	89

ABSTRACT

State space models are widely used for conducting time series analysis. Developing a state space model involves proposing mathematical equations that describe how a data-generating system evolves over time and how observations of the system are obtained. These models are particularly useful when a scientific hypothesis about system dynamics exists, as is common when modeling ecological populations or tracking infectious disease outbreaks over time. However, except for the simplest cases, state space models do not permit closed-form expressions of their likelihood functions, presenting challenges for inference. This thesis presents three projects that introduce innovations in likelihood-based inference for state space models.

The first project proposes a novel approach for performing inference on Auto Regressive Moving Average (ARMA) time series models, which are formally linearly Gaussian state space models. ARMA models are among the most frequently taught and widely used methods for time series analysis. In this project, I demonstrate that existing algorithms and software for parameter estimation often produce sub-optimal parameter estimates with surprising frequency. I introduce a novel random initialization algorithm designed to leverage the structure of the ARMA likelihood function to help overcome these optimization shortcomings. Additionally, I demonstrate that profile likelihoods offer superior confidence intervals compared to those based on the Fisher information matrix—the current standard practice for ARMA modeling.

The second project presents a likelihood-based analysis of the 2010-2019 cholera outbreak in Haiti. This work explores three distinct state space models for cholera incidence data and demonstrates the effectiveness of recently developed algorithms for performing inference in a high-dimensional setting. A key focus of this project is to assess the strengths and limitations of using state space models to inform public health policy decisions. Existing methodologies and workflows for this purpose are evaluated, and revised data analysis strategies that lead to better statistical fit and outcomes are presented. For example, I demonstrate a reproducible framework for diagnosing model misspecification and subsequently developing enhancements that result in better recommendations for policy decisions.

The third project proposes a simulation-based algorithm designed to perform maximum

likelihood estimation for a class of high-dimensional state space models. This algorithm, called the Marginalized Panel Iterated Filter (MPIF), significantly enhances the capability of iterated filtering algorithms to estimate parameters for large collections of independent but related state space models. Improvements in parameter estimates and empirical convergence rates are achieved by addressing the issue of particle depletion that occurs when performing iterated filtering on models that have high-dimensional parameter spaces. Theoretical support for the algorithm is provided through an analysis of iterating marginalized Bayes maps. Additionally, asymptotic theory demonstrating the convergence of general iterated filtering algorithms for panel models without the marginalization step is presented.

CHAPTER 1

Introduction

The importance of time series analysis has continued to grow with the demand for both the collection and analysis of data. Time series data are abundant, as data are often generated over time in a way that introduces dependence between observations. However, the analysis of time series presents a unique challenge for both statisticians and practitioners, as the temporal dependence between observations requires careful treatment. A unifying approach in time series analysis, known as state space modeling, assumes that observations from the system under investigation depend on an unobservable process. Constructing a state space model involves proposing equations that describe the evolution of this unobservable process over time and how observations are made. Common goals of state space modeling include making inferences about critical yet unobserved variables, incorporating real-world mechanisms into statistical models, and enabling likelihood based inference for models who that lack closed-form expressions of the likelihood function.

The term state space model (SSM) has a number of synonyms including mechanistic model, partially observed Markov process (POMP) model, and hidden Markov model (HMM) (King et al., 2015). In some literature, these terms have alternatively been used to refer to particular instances of the more general model class (e.g., Glennie et al., 2023), resulting in some confusion. Throughout the dissertation, I follow the seminal work of Durbin and Koopman (2012) and use the term SSM to refer to the most general case. That is, the observed time series is assumed to be reliant on an unobservable collection of variables at each observation time, without any restrictions on relationships between variables. The term *mechanistic* is reserved for models where the evolution of latent variables is dictated by equations intended to replicate real-world mechanisms, such as the transmission of infectious diseases among susceptible populations (Wheeler et al., 2024). In most instances in this thesis and existing literature, the latent variables are treated as a Markov process, in which case the term POMP can be used. Finally, existing literature often uses the term HMM to refer to POMP models where the latent variables take values in a discrete and finite space (Eddy,

2004; Doucet et al., 2001; Glennie et al., 2023; Newman et al., 2023); this same definition is used here.

1.1 Defining State Space Models and Notation

In this section, some general terminology and notation that is common across chapters is introduced, and refined in subsequent chapters as needed. This thesis presents methodologies for modeling real-valued time series data, collected at observation times $t_1, t_2, \dots, t_N \in \mathcal{T} \subset \mathbb{R}$. Observation times do not need to be equally spaced. Notably, this treatment includes time series where measurements are equally spaced but occasionally missing.

The observation at time t_n , $n \in \{1, \dots, N\}$ can be vector valued, and is denoted $\mathbf{y}_n^* \in \mathbb{R}^{d_y}$. The collection of observations $\{\mathbf{y}_n^*, n \in 1, \dots, N\}$ is modeled as a single realization of a collection of random variables $\{\mathbf{Y}_n, n \in 1, \dots, N\}$, called the observable or measurement process. The state space formulation introduces a collection of unobservable random variables $\{\mathbf{X}_n, n \in 1, \dots, N\}$ that exist at the same time points as the observable process, and we generally assume that this unobservable (or latent) process takes values in a subset of \mathbb{R}^{d_x} . Most often, models of interest include a description of how the latent process is initialized at some time $t_0 < t_1$, and the collection of latent variables is extended to include \mathbf{X}_0 . The latent process can be modeled either as a discrete or continuous time process.

For convinience, we adopt the notation that for any integers $a < b$, $a:b$ is the vector $(a, a+1, \dots, b-1, b)$. Similarly, we write the entire collection of observations as $\mathbf{y}_{1:N}^* = (\mathbf{y}_1^*, \dots, \mathbf{y}_N^*)$, and use the same basic notation for $\mathbf{Y}_{1:N}$ and $\mathbf{X}_{0:N}$. The joint probabiliy density or (mass function) of the observable and unobservable processes is assumed to exist, and can be written as $f_{\mathbf{X}_{0:N}, \mathbf{Y}_{1:N}}(\mathbf{x}_{0:N}, \mathbf{y}_{1:N}; \theta)$, where θ is a parameter vector $\theta \in \mathbb{R}^{d_\theta}$. The likelihood is a function of θ , defined by the density of $\mathbf{Y}_{1:N}$ evaluated at the observed data $\mathbf{y}_{1:N}^*$:

$$\mathcal{L}(\theta) = f_{\mathbf{Y}_{1:N}}(\mathbf{y}_{1:n}^*; \theta) = \int f_{\mathbf{X}_{0:N}, \mathbf{Y}_{1:N}}(\mathbf{x}_{0:N}, \mathbf{y}_{1:N}^*; \theta) d\mathbf{x}_{0:N}.$$

In all but the simplest models, a closed-form expression of the likelihood function is not readily available due to the high-dimensional integral involving the latent variables.

The state space models described in Chapters 2–?? can all be classified as POMP models. This model class makes some additional, unrestrictive assumptions that can be used to enable likelihood based inference. The latent process is assumed to be Markovian and independent of the measurement process. Additionally, the measurements are assumed to be conditionally independent given the current value of the latent process. In terms of joint and conditional density functions of variables at observation times, these assumptions imply that for all

$n \in 1:N$,

$$f_{\mathbf{x}_n|\mathbf{x}_{1:n-1}}(\mathbf{x}_n|\mathbf{x}_{1:n-1};\theta) = f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1};\theta),$$

and

$$f_{\mathbf{Y}_n|\mathbf{x}_{1:N}, \mathbf{Y}_{1:n-1}, \mathbf{Y}_{n-1:N}}(\mathbf{y}_n|\mathbf{y}_{1:N}, \mathbf{y}_{1:n-1}, \mathbf{y}_{n-1:N};\theta) = f_{\mathbf{Y}_n|\mathbf{x}_n}(\mathbf{y}_n|\mathbf{x}_n;\theta).$$

TODO: Add a new figure here representing POMP models.

Though these additional assumptions do not generally give rise to a closed-form expression of the likelihood function, they do enable likelihood based inference in many scenarios. For instance, if the evolution of the latent variables is dictated by linear Gaussian equations, then the Kalman filter (Kalman, 1960) can be used to evaluate the likelihood function. More generally, these assumptions enables likelihood based inference via *plug-and-play* algorithms, which only require the ability to simulate the latent process (Bretó et al., 2009). The additional assumptions of a POMP model can also be used to refactor the joint density of latent and observable random variables as

$$f_{\mathbf{x}_{0:N}, \mathbf{Y}_{1:N}}(\mathbf{x}_{0:N}, \mathbf{y}_{1:N};\theta) = f_{\mathbf{x}_0}(\mathbf{x}_0;\theta) \prod_{n=1}^N f_{\mathbf{Y}_n|\mathbf{x}_n}(\mathbf{y}_n|\mathbf{x}_n;\theta) f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1};\theta).$$

This representation is useful as the entire model can be expressed in terms of three simple components, namely the *initializer* $f_{\mathbf{x}_0}(\mathbf{x}_0;\theta)$, the *process* or *transition* model $f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1};\theta)$, and the *measurement* model $f_{\mathbf{Y}_n|\mathbf{x}_n}(\mathbf{y}_n|\mathbf{x}_n;\theta)$. Each of these pieces can depend arbitrarily on observation time and the individual components of the parameter vector θ .

1.2 Overview of remaining chapters

The remainder of the thesis is comprised of three distinct chapters that have been published in peer-reviewed articles or have been prepared to be published. Each chapter contains includes some work from coauthors of these papers; major contributions from coauthors are highlighted in each chapter they appear.

Chapter 2 describes an innovation to the maximum likelihood estimation procedure of autoregressive moving average (ARMA) models. ARMA models are among the most frequently used methods for analyzing time series data. Although they are not typically considered state space models, the parameters to ARMA models are often estimated by using an equivalent linear Gaussian state space representation of the model. In this chapter, I provide a summary of how parameter estimates are obtained via maximum likelihood and discuss common pitfalls that may lead to sub-optimal parameter values. Using various simulation studies, I demonstrate that a newly proposed random restart algorithm may lead

to higher model likelihoods.

In Chapter 3, I analyze the 2010-2019 cholera outbreak in Haiti using three POMP models. Public health decisions must be made about when and how to implement interventions to control an infectious disease epidemic. These decisions should be informed by data on the epidemic as well as current understanding about the transmission dynamics. Such decisions can be posed as statistical questions about scientifically motivated dynamic models. Existing standards for time series modeling of infectious diseases are insufficient to consistently make reliable conclusions using these models. In this chapter, I summarize current approaches of fitting and evaluating dynamic models, and develop data analysis strategies that lead to improved statistical fit. Specifically, I present approaches for the diagnosis of model misspecification, development of alternative models, and computational improvements in optimization, in the context of likelihood-based inference on nonlinear dynamic systems. The work contained in this chapter has been submitted to the *Annals of Applied statistics*, and is currently under review (Wheeler and Ionides, 2023).

In the final chapter, I propose a new algorithm called the marginalized panel iterated filter (MPIF) that can be used to calibrate model parameters for PanelPOMP models. The MPIF algorithm takes advantage of the special structure of independence in PanelPOMP models to outperform other algorithms that can be used in a more general setting. The algorithm permits the distinction between shared and unit-specific parameters, which is necessary in many scientific applications. While theoretical arguments demonstrating that the algorithm converges to the maximum likelihood estimate have yet to be derived, comparisons against algorithms that do have theoretical guarantees suggest that MPIF may indeed maximize model likelihoods.

CHAPTER 2

Likelihood Based Inference for ARMA Models

2.1 Introduction

Auto-regressive moving average (ARMA) models are the most well known and frequently used approach to modeling time series data. The general ARMA model was first described by Whittle (1951), and popularized by Box and Jenkins (1970). Today, ARMA models are a fundamental component of various academic curricula, leading to their widespread use in both academia and industry. ARMA models are as foundational to time series analysis as linear models are to regression analysis, and they are often used in conjunction for regression with ARMA errors. A Google Scholar search for articles from 2024 onward that include the phrase “time series” and the term “ARMA” (or variants) yields over 18,000 results. While not all these articles focus on the same models discussed here, the importance of this model class to modern science cannot be overstated. Given the ubiquity of ARMA models, even small improvements in parameter estimation constitute a significant advancement of statistical practice.

A commonly used extension of the ARMA model is the *integrated* ARMA model, which extends the class of ARMA models to include first or higher order differences. That is, an autoregressive integrated moving average (ARIMA) model is an ARMA model fit after differencing the data in order to make the data stationary. Additional extensions include the modeling of seasonal components (SARIMA), or the inclusion of external regressors (SARIMAX). Our methodology can readily be extended to these model classes as well, but here we focus on ARMA modeling for simplicity.

We demonstrate that the most commonly used methodologies and software for estimating ARMA model parameters frequently yield sub-optimal estimates. This assertion may seem surprising given the extensive application and study of ARMA models over the past five decades. A natural question arises: if these optimization issues exist, why have they not been addressed? There are three plausible explanations: the first is that the potential for

sub-optimal results has largely gone unnoticed; the second is a satisfactory solution has not yet been discovered; and the third is a general indifference to the problem. It is likely that a combination of these factors has deterred prior exploration of this issue. For instance, most practitioners may be unaware of the problem, while those who have noticed it either did not prioritize it or were unable to provide a general method to resolve it. In this article, we address all three possible explanations by demonstrating the existence of an existing shortcoming, explaining why the problem has nontrivial consequences, and proposing a readily applicable and computationally efficient solution.

Imperfect likelihood optimization has an immediate consequence of complicating standard model selection procedures. Algorithms in widespread use lead to frequent inconsistencies in which a smaller model is found to have a higher maximized likelihood than a larger model within which it is nested. This is mathematically impossible but occurs in practice when the likelihood is imperfectly maximized, and is commonly observed using contemporary methods for ARMA models. Such inconsistencies are a distraction for the interpretation of results even when they do not substantially change the conclusion of the data analysis. Removing numeric inconsistencies can, and should, increase confidence in the correctness of statistical inferences.

There are various software implementations available for the estimation of ARMA model parameters. In this article, we focus on the standard implementations in R (`stats` package) and Python (`statsmodels` module), which we selected due to their widespread usage. While both implementations offer multiple ways to estimate parameters, the default approach in both software packages is to perform likelihood maximization, assuming that the error terms are Gaussian. The challenges in parameter estimation arise in this situation because there is no closed-form expression for the likelihood function, though computational algorithms do exist for maximizing ARMA likelihoods Gardner et al. (1980).

We begin by providing essential background information on the estimation of ARMA model parameters. We then present our proposed approach for parameter estimation, which leads to parameter values with a likelihood that is never lower and sometimes higher than the standard method. This is followed by a motivating example and a discussion of the potential implications of our proposed method. We also discuss the construction of standard errors for our maximum likelihood estimate. Specifically, we show that estimates of the standard error for model parameters that are default output of R and Python can be misleading, and we provide a reliable alternative. Throughout the article, we use the `stats::arima` function from the R programming language as a baseline for comparison. The same methodology for fitting parameters is used in the `statsmodels.tsa` module in Python, and we demonstrate that our results apply to this software as well (Appendix A.1).

2.2 Maximum Likelihood for ARMA Models

Following the notation of Shumway and Stoffer Shumway and Stoffer (2017), a time series $\{x_t; t = 0, \pm 1, \pm 2, \dots\}$ is said to be ARMA(p, q) if it is (weakly) stationary and

$$x_t = \phi_1 x_{t-1} + \dots + \phi_p x_{t-p} + w_t + \theta_1 w_{t-1} + \dots + \theta_q w_{t-q}, \quad (2.1)$$

with $\{w_t; t = 0, \pm 1, \pm 2, \dots\}$ denoting a mean zero white noise (WN) processes with variance $\sigma_w^2 > 0$, and $\phi_p \neq 0, \theta_q \neq 0$. We refer to the positive integers p and q of Eq. 2.1 as the autoregressive (AR) and moving average (MA) orders, respectively. A non-zero intercept could also be added to Eq. 2.1, but for simplicity we assume that the time series $\{x_t; t = 0, \pm 1, \pm 2, \dots\}$ has zero mean. We denote the set of all model parameters as $\psi = \{\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q, \sigma_w^2\}$. Our objective is to estimate model parameters using the observed uni-variate time series data.

Given the importance of ARMA models, numerous methods have been developed for parameter estimation. For instance, parameters can be estimated through methods such as Bayesian inference Monahan (1983); Chib and Greenberg (1994) or neural networks Chon and Cohen (1997), among many others. Specialized methods include those for integer-valued data Weiß and Zhu (2024), approaches robust to outliers Chakhchoukh (2010), and non-Gaussian ARMA processes Lii (1990). While each methodology has its merits, we focus on the approach most widely adopted in statistical practice: likelihood maximization, under the assumption that the WN process $\{w_t\}$ is Gaussian.

A relevant method for estimating model parameters is minimization of the conditional sum-of-squares (CSS). The CSS estimate is fast to compute, but it does not possess the statistical efficiency of the maximum likelihood estimate (MLE). However, the CSS method plays a role in likelihood maximization, so we briefly describe it here. By solving for the WN term w_t , Eq. 2.1 can be written as

$$w_t = x_t - \sum_{i=1}^p \phi_i x_{t-i} - \sum_{j=1}^q \theta_j w_{t-j}. \quad (2.2)$$

A natural estimator would involve minimizing the sum of squares $\sum_{t=1}^n w_t^2$. However, since only x_1, x_2, \dots, x_n are observed and w_t is recursively defined in Eq. 2.2 using values of x_{t-p} and w_{t-q} , directly minimizing this sum is intractable. The CSS method addresses this issue by conditioning on the first p values of the process, assuming $w_p = w_{p-1} = \dots = w_{p+1-q} = 0$, and minimizing the conditioned sum $\sum_{t=p+1}^n w_t^2$. While the CSS method provides an attractive solution due to its relative simplicity and easiness to compute, it ignores the error terms for the

first few observations. This is particularly concerning when the time series is short or when there are missing observations. CSS minimization was previously popular because methods for likelihood maximization were considered prohibitively slow, though this is no longer the case with currently available hardware and software Ripley (2002).

For likelihood maximization, Eq. 2.1 is reformulated as an equivalent state-space model. Although there are several ways this can be done, the approach of Gardner et al. (1980) is widely used. In this approach, we let $r = \max(p, q + 1)$ and extend the set of parameters so that $\bar{\psi} = \{\phi_1, \dots, \phi_r, \theta_1, \dots, \theta_{r-1}, \sigma_w^2\}$, with some of the ϕ_i s or θ_i s being equal to zero unless $p = q + 1$. We define a latent state vector $z_t \in \mathbb{R}^r$, along with transition matrices $T \in \mathbb{R}^{r \times r}$ and $Q \in \mathbb{R}^{r \times 1}$, enabling the recovery of the original sequence $\{x_t\}$ using Equations 2.3 and 2.4. For a detailed explanation of how to define the latent state z_t and transition matrices T and Q in order to recover the ARMA model, we refer readers to Chapter 3 of Durbin and Koopman Durbin and Koopman (2012). Along with initializations for the mean and variance of z_0 , these equations allow for the exact computation of the likelihood of the ARMA model via the Kalman filter Kalman (1960), which can subsequently be optimized by a numeric procedure such as the BFGS algorithm Fletcher (2000).

$$z_t = Tz_{t-1} + Qw_t, \quad (2.3)$$

$$x_t = \begin{pmatrix} 1 & 0 & \dots & 0 \end{pmatrix} z_t. \quad (2.4)$$

Numeric black-box optimizers require an initial guess for parameter values. For ARMA models, this task is non-trivial because the valid parameter region is defined in terms of the roots of a polynomial associated with the parameters, as discussed in the next section. The default strategy in R and Python is to use the CSS estimator for initialization. This is an effective approach because the CSS estimator asymptotically converges to the MLE Shumway and Stoffer (2017), and may therefore be close to the global maximum when there are sufficiently many observations. However, the CSS initialization is less useful with limited data, or when there are missing observations. The CSS estimate may also lie outside the valid parameter region, and in such cases, parameters are reinitialized at the origin. Both software implementations also allow for manual selection of initial parameter values, but finding suitable initializations manually can be challenging due to complex parameter inter-dependencies.

The log-likelihood function of ARMA models is often multimodal Ripley (2002), and therefore this single initialization approach can result in parameter estimates corresponding to local maxima (see Fig 2.1). This is true even for a carefully chosen initialization, such as the CSS estimate. A common strategy to optimize multimodal loss functions is to

perform multiple optimizations using different initial parameters. However, we have found no instances of practitioners using a multiple initialization strategy for estimating ARMA model parameters. This may be explained by a general unawareness of the possibility of converging to a local maximum or because obtaining a suitable collection of initializations for ARMA models is nontrivial. For example, independently initializing parameters at random can place the parameter vector outside the region of interest. Furthermore, uniform random sampling generally fails to adequately cover the plausible parameter region (Appendix A.2).

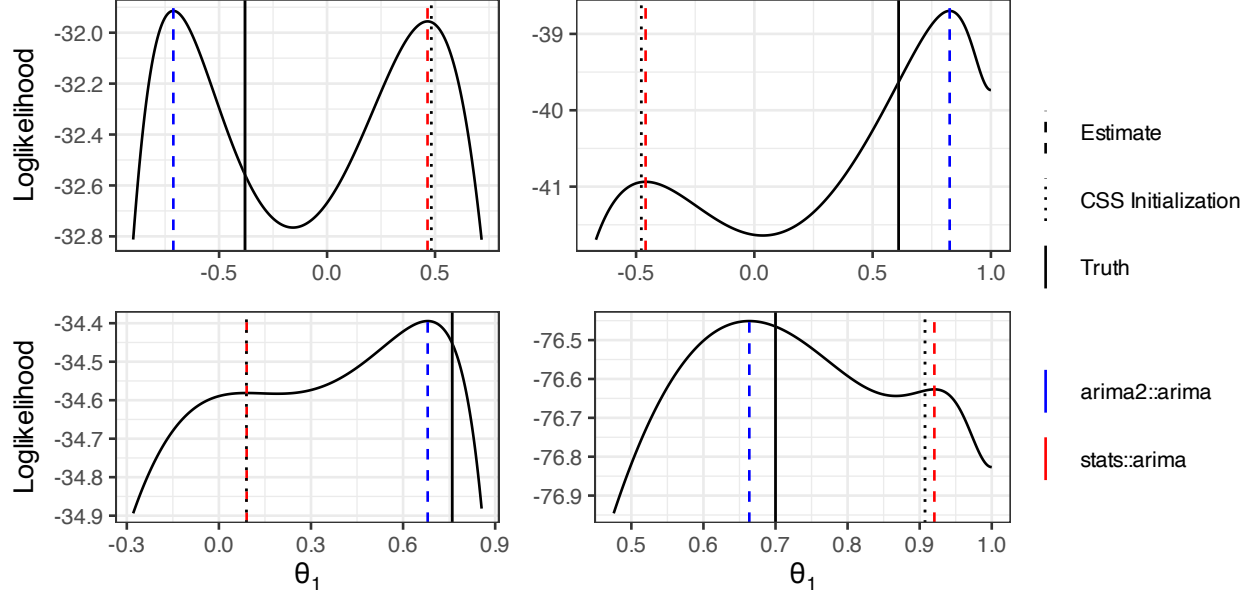


Figure 2.1: The profile log-likelihood of data simulated from four distinct MA(1) models, demonstrating a few examples of multimodal likelihood surfaces. The solid, black line indicates the true value of θ_1 ; the dotted line is the CSS-initialization. The dashed lines correspond to the estimate $\hat{\theta}_1$ using `stats::arima` (red) and our proposed algorithm (implemented in `arima2::arima`, blue).

2.2.1 A Novel Multi-start Algorithm

To obtain random parameter initializations, parameter sets must correspond to *causal* and *invertible* ARMA processes; definitions are in Chapter 3 of Shumway and Stoffer Shumway and Stoffer (2017). Let $\{\phi_i\}_{i=1}^p$ and $\{\theta_i\}_{i=1}^q$ be the coefficients of the ARMA(p, q) model (Eq. 2.1), and define $\Phi(x) = 1 - \phi_1 x - \phi_2 x^2 - \dots - \phi_p x^p$ as the AR polynomial, and $\Theta(x) = 1 + \theta_1 x + \theta_2 x^2 + \dots + \theta_q x^q$ as the MA polynomial. An ARMA model is *causal* and *invertible* if the roots of the AR and MA polynomials lie outside the complex unit circle. Therefore in order to obtain valid, random parameter initializations, we sample the roots of

$\Phi(x)$ and $\Theta(x)$ and use these values to reconstruct parameter initializations.

It is an easier task to sample *inverted* roots, as the sufficient conditions for causality and invertibility now require that the inverted roots lie inside the complex unit circle, a region easier to sample uniformly. The roots of the polynomials can be real or complex; complex roots must come in complex conjugate pairs in order for all of the corresponding model parameters to be real. The simplest approach would be to sample all inverted root pairs (z_1, z_2) within the complex unit circle by uniformly sampling angles and radii (lines 10-14 of Algorithm 1). However, this would imply almost surely all root pairs are complex, and some model parameters would only be sampled as positive (or negative). For instance, consider an AR(2) model. The AR polynomial is:

$$\Phi(x) = 1 - \phi_1 x - \phi_2 x^2 = (1 - z_1 x)(1 - z_2 x) = 1 - (z_1 + z_2)x - (z_1 z_2)x^2$$

In this equation, if both z_1, z_2 are complex conjugates, then $\phi_2 = z_1 z_2 > 0$. As such, the only way that $\phi_2 < 0$ is if $z_1, z_2 \in \mathbb{R}$. Similar results hold for the MA coefficients with opposite signs for the coefficients. This issue is directly addressed in lines 5-8 of Algorithm 1: root pairs are sampled as real with probability $p = \sqrt{1/2}$, and real pairs sampled with the same sign with probability p , such that the product (and sums) of each pair is positive with probability $1/2$. We sample conjugate pairs within an annular disk on the complex plane to avoid trivial and approximately non-stationary cases. The radii of both the inner and outer circles defining the disk are defined using γ in lines 7, 10, and 13 of Algorithm 1.

Parameter redundancy in ARMA models occurs when the polynomials $\Phi(x)$ and $\Theta(x)$ share one or more roots, leading to an overall reduction in model order. This complicates parameter initialization, optimization, and identifiability. The ARMA model (Eq. 2.1) can be rewritten as:

$$\Phi(B)x_t = \Theta(B)w_t, \tag{2.5}$$

where B is the *backshift* operator, i.e., $Bx_t = x_{t-1}$. Using the fundamental theorem of algebra, Equation 2.5 can be factored into

$$(1 - \lambda_1 B) \dots (1 - \lambda_p B)x_t = (1 - \nu_1 B) \dots (1 - \nu_q B)w_t,$$

where $\{\lambda_i\}_{i=1}^p$ and $\{\nu_j\}_{j=1}^q$ are the inverted roots of $\Phi(B)$ and $\Theta(B)$, respectively. If $\lambda_i = \nu_j$ for any $(i, j) \in \{1, \dots, p\} \times \{1, \dots, q\}$, then the roots will cancel each other out, resulting in an ARMA model of smaller order. As an elementary example, consider the ARMA(1, 1) and

Algorithm 1: MLE for ARMA Models.**Inputs (defaults):**First parameter initialization $\psi_0 = (\phi_1^0, \dots, \phi_p^0, \theta_1^0, \dots, \theta_q^0)$ (CSS estimate).Minimum acceptable polynomial root distance $\alpha > 0$, ($\alpha = 0.01$).Probability of sampling a root pairs as real $0 \leq p \leq 1$, ($p = \sqrt{1/2}$).Bounds on inverted polynomial roots $\gamma \in (0, 0.5)$, ($\gamma = 0.05$).Numeric optimization routine $f(\psi)$ Gardner et al. (1980).Stopping Criterion (stop if last M iterations do not improve log-likelihood $\ell(\psi)$).

```

1  Get preliminary estimate:  $\hat{\psi}_0 = f(\psi_0)$ ; set  $k = 0$ ;
2  repeat Until stopping criterion met
3      Set AR and MA roots  $\{z_i^{\text{AR}}\}_{i=1}^p = 0_p$ ,  $\{z_i^{\text{MA}}\}_{i=1}^q = 0_q$ ; increment  $k$ ;
4      while  $\min_{i,j} |z_i^{\text{AR}} - z_j^{\text{MA}}| < \alpha$ , for both AR and MA polynomials do
5          Sample paired roots as real with probability  $p$ ;
6          for all real pairs do
7              Sample root magnitudes from  $U(\gamma, 1 - \gamma)$ ;
8              Sample signs with  $P(\text{sign}(z_1) = \text{sign}(z_2)) = p$ ;
9          for all complex pairs do
10             sample angle:  $\tau \sim U(0, \pi)$ ; sample radius:  $r \sim U(\gamma, 1 - \gamma)$ ;
11             set  $z_1 = r \cos(\tau) + ir \sin(\tau)$ ; set  $z_2 = \bar{z}_1$ ;
12             if Number of roots is odd (non-paired root) then
13                 sample  $\tau$  uniformly from the set  $\{0, \pi\}$ ; sample  $r \sim U(\gamma, 1 - \gamma)$ ;
14                 set  $z = r \cos(\tau)$ ;
15             Calculate coefficients  $\psi_k = (\phi_1^k, \dots, \phi_p^k, \theta_1^k, \dots, \theta_q^k)$  using sampled roots;
16             Estimate  $\hat{\psi}_k = f(\psi_k)$ ;
17 until;
18 Set  $\hat{\psi} = \arg \max_{j \in 0:k} \ell(\hat{\psi}_j)$ ;

```

ARMA(2, 2) models in equations 2.6 and 2.7.

$$x_t = \frac{1}{3}x_{t-1} + w_t + \frac{2}{3}w_{t-1}, \quad (2.6)$$

$$x_t = \frac{5}{6}x_{t-1} - \frac{1}{6}x_{t-2} + w_t + \frac{1}{6}w_{t-1} - \frac{1}{3}w_{t-2}. \quad (2.7)$$

While these two models appear distinct at first glance, re-writing the models in polynomial form (Eq. 2.5) shows that these two models are actually equivalent after canceling out the common factors on each side of the equation.

In a similar fashion, it is possible that the roots are not exactly equal but are approximately equal. In this case, the ratio of factors becomes close to one, resulting in a similar effect to when the roots exactly cancel. We avoid the possibility of *nearly canceling roots* in parameter initializations by requiring the minimum Euclidean distance between inverted polynomial

roots to be greater than α . This is done in line 4 of Algorithm 1, though the condition is rarely triggered if the order of the model is of typical size ($p, q < 4$).

Our sampling scheme is combined with existing procedures for numeric optimization of model log-likelihoods (lines 1 and 16), as well as the default initialization strategy for ψ_0 used by existing software (such as the CSS initialization). Doing so guarantees that final estimates correspond to likelihood values greater than or equal to the currently accepted standards in the software environment where the algorithm is implemented. For this article, both the numeric optimization procedure $f(\cdot)$ and the parameter initialization strategy to obtain ψ_0 are those implemented in `stats::arima`. The stopping criterion was chosen so that the algorithm stops trying new initial values when no new maximum has been found using the last M parameter initializations. Alternative stopping criterion can be used (see for example, Hart (1998)), but we found that this simple heuristic works well in practice.

Algorithm 1 is implemented in the R package `arima2`, available on the Comprehensive R Archive Network (CRAN) Wheeler et al. (2023). The package features the function `arima2::arima`, which is an adaptation of the `stats::arima` function modified to incorporate the adjustments specified by Algorithm 1.

2.2.2 Simulation Studies

To investigate the extent to which the standard approach for ARMA parameter estimation results in improperly maximized likelihoods, we conduct a series of simulation studies. It is challenging to obtain precise estimates of how frequently current standards lead to sub-optimal parameter estimates due to the varied applications of ARMA models in practice, the diversity in data sizes (n) and model orders (p, q), and the differing degrees to which an ARMA model adequately describes the data-generating process. Therefore, we restrict our simulation studies to idealized scenarios where the data-generating process is Gaussian-ARMA, recognizing that likelihood maximization is easiest for this model class, thereby resulting in conservative estimates of how frequently our algorithm improves model likelihood.

In the first simulation study, we simulate time series data of lengths $n \in \{50, 100, 500, 1000\}$ from Gaussian-ARMA models with known orders $(p, q) \in \{1, 2, 3\}^2$, generating 36,000 unique models and datasets. We avoid any models that contain parameter redundancies by requiring the data generating model to have a minimum distance of 0.1 between all roots of $\Phi(x)$ and $\Theta(x)$. We further restrict model coefficients so that they do not lie near boundary conditions. Models of the same order of the generating data are fit to the data, simplifying the the problem further by avoiding the order selection step that is necessary in most data analyses. In doing so, we attempt to answer the question of how often sub-optimal estimates may arise

using existing software in the case where the parameter estimation procedure should be as easy as possible for the given combinations of (n, p, q) .

Even in this extremely simplified scenario, existing software failed to properly maximize model likelihoods in at least 20.8% of the simulated datasets—evidenced by an improvement obtained using Algorithm 1. Though this improvement may appear modest, an improvement in 20.8% of the large number of published ARMA models would affect many papers—a number measured in thousands of papers since 2024 following our estimate in the introduction. Furthermore, time series analysis courses and textbooks often recommend fitting multiple ARMA models to a dataset, and here we only fit one for each algorithm. Consequently, the probability that at least one candidate model is not properly optimized increases significantly in practice. The rate of improvement obtained using our algorithm increases with model complexity and decreases with more observations (Fig. 2.2). For example, likelihoods improved in 55.1% of the simulations when $p = q = 3$ and $n = 50$; models of this size and number of observations are not uncommon in published research studies.

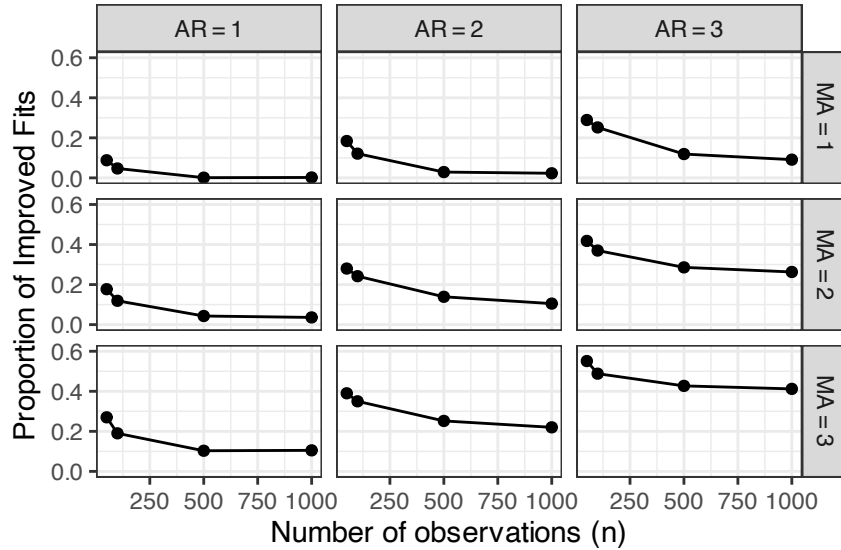


Figure 2.2: Proportion of simulated data with improved likelihood from using multiple restarts (Algorithm 1).

Importantly, we do not claim to improve likelihood for all ARMA models and datasets. In this simulation study, our likelihood maximum routine did not improve likelihoods for many of the simulated data. However, these results demonstrate that there is a very real potential for obtaining sub-optimal parameter estimates when using only a single parameter initialization, even in the most idealistic scenarios. Rather than having to worry if a single initialization is sufficient to fit a given model, it is preferable to adopt methods that make such

situations rare. The primary limitation of our algorithm is that the potential for improved fits comes at the cost of increased computation times. In our simulation study, however, the average time to estimate parameters using our approach was 0.6 seconds, a computational expense that is worth the effort in many situations.

The median log-likelihood improvement in this simulation study was 0.66, with an interquartile range of (0.22, 1.46). Among the most common motivations for fitting ARMA models is to model serial correlations in a regression model; in this setting, the discovered shortcomings in log-likelihood are often enough to change the outcome of the analysis. For instance, consider modeling $y_i = \beta x_i + \epsilon_i$, where $\beta \in \mathbb{R}$, and we model the error terms $\epsilon_i \sim \text{ARMA}(p, q)$. For now, we will assume the order (p, q) is fixed. We may wish to test the hypothesis $H_0 : \beta = 0$ vs $H_1 : \beta \neq 0$. A standard approach to doing this is a likelihood ratio test, and using Wilks' theorem to get an approximate test. We denote ll_0 and ll_1 as the maximum log-likelihood of the model under H_0 and H_1 , respectively. The standard approximation is to assume $2\Delta = 2(ll_1 - ll_0) \sim \chi_1^2$. Using a significance level of $\alpha = 0.05$, we would reject H_0 if $\Delta \geq 1.92$. Given that $E_{H_0}[\Delta] = 0.5$, subtracting the reported log-likelihood deficiencies (which has a median value of 0.66) of existing software to either or both ll_0, ll_1 could change the outcome of this test.

2.2.2.1 Parameter Uncertainty

Improved parameter estimation leads to modified standard error estimates, which are default outputs in R and Python. These standard errors result from the numeric optimizer's estimate of the gradient of the log-likelihood, used to approximate the Fisher information matrix. These standard errors, though not inherently of interest, are sometimes used justify the inclusion of a parameter in a model (Brockwell and Davis, 1991, Chapter 9). We extend our simulation study to examine this approach and how our algorithm impacts the estimates. For each of the 36,000 generative models from the previous study, we generate 100 additional datasets, estimating the MLE and Bonferroni-adjusted 95% confidence intervals using both estimated standard errors and profile likelihood confidence intervals (PLCIs) from Wilks' theorem. Fig 2.3 shows PLCIs had better or equivalent nominal coverage than Fisher-based confidence intervals across all combinations of p , q , and n . We found that the interval estimation method mattered more for confidence interval performance than the specific parameter estimation algorithm. The relevance of this result is explored further in Section 2.3.1.

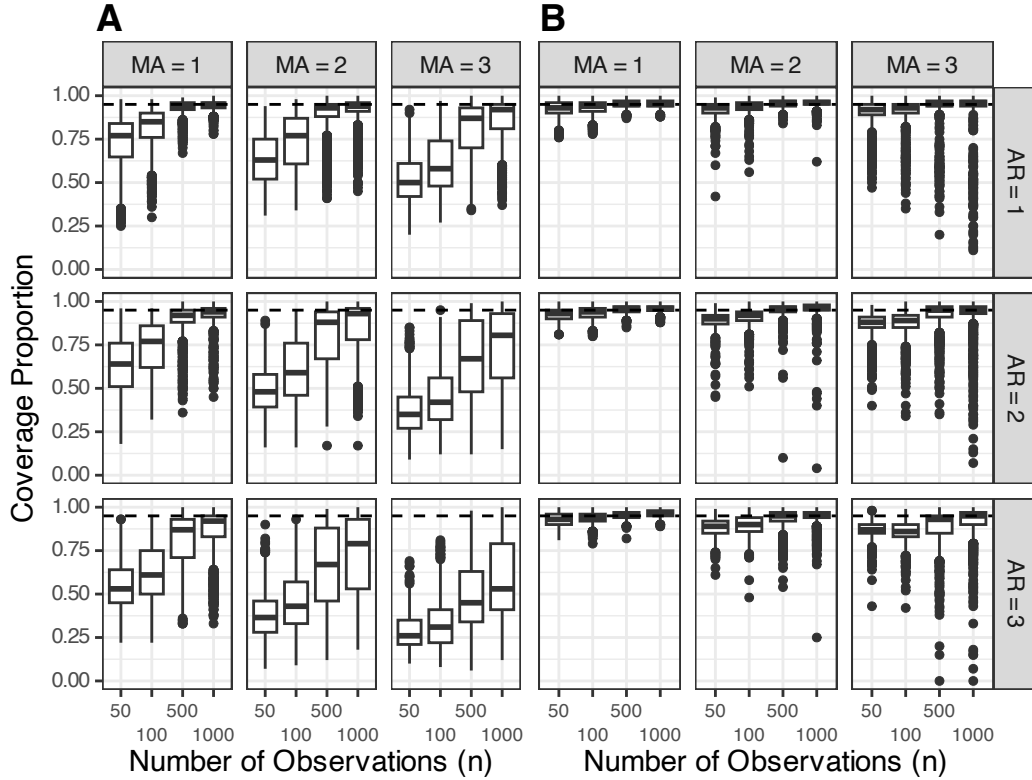


Figure 2.3: Proportion of models that achieved nominal coverage of Bonferroni adjusted 95% confidence intervals. The dashed line denotes the target coverage level. (A) Confidence intervals created using Fisher’s information matrix. (B) Confidence intervals created using profile likelihoods.

2.2.2.2 AIC Table Consistency

A more realistic situation than the previous simulation study involves estimating the model order (p, q) as well as obtaining parameter estimates. Fitting multiple models raises the chance that at least one candidate model was not properly optimized. It may also necessitate fitting larger models than needed, leading to parameter redundancies that make proper optimization more challenging.

A contemporary approach involves fitting several candidate models and selecting the one that minimizes a criterion like Akaike’s information criterion (AIC) Akaike (1974). This can be done by explicitly creating a table of all candidate models and their corresponding AIC values; in this case issues of improper maximization become more apparent. For instance, a table of AIC values may contain numeric inconsistencies, where a larger model may have lower estimated likelihoods than a smaller model within which it is nested (for an example, see Section 2.3). This type of result can make a careful practitioner feel uneasy, as there is

evidence that at least one candidate model was not properly optimized. Evidence of improper optimization may be less evident when relying on software that automates this process, such as the automated Hyndman-Khandakar algorithm Hyndman and Khandakar (2008), but the potential for sub-optimal estimates remains.

We conducted an additional simulation study to investigate numeric inconsistencies that may arise when fitting multiple model parameters. As before, we simulated 1000 unique models and datasets of size $n \in \{50, 100, 500, 1000\}$ from Gaussian ARMA(p, q) models for $(p, q) \in \{1, 2, 3\}^2$. To avoid models with parameter redundancies, we ensured a minimum distance of 0.1 between all roots of $\Phi(x)$ and $\Theta(x)$ and excluded models with coefficients near boundary conditions. For each dataset, AIC tables were created for model sizes $(p, q) \in \{0, 1, 2, 3\}^2$.

The single parameter initialization approach resulted in AIC table inconsistencies in 45.6% of the simulated datasets. Although our proposed algorithm significantly mitigates this issue, it does not guarantee that all model likelihoods are fully maximized. This is illustrated in Fig 2.4, where a non-zero percentage of AIC tables remain inconsistent, even as the algorithm’s stopping criterion grows. The ARMA(1, 1) panel in Fig 2.4 illustrates the increasing difficulty of parameter estimation when dealing with parameter redundancies. In such cases, it is often necessary to adjust additional parameters in the numeric optimization routine. For example, our R implementation of the algorithm relies on the generic BFGS optimizer in the `stats::optim` function. Modifying the optimization method or the default hyperparameters can lead to improved fits or faster convergence rates. While the default parameters of the numeric optimizer are generally adequate, increasing the maximum number of algorithmic iterations can be beneficial for fully maximizing the likelihood for challenging models and data.

The contemporary approach of using AIC—or any other information based criteria—to select model order involves fitting unnecessarily large models, leading to parameter redundancies that complicate likelihood optimization. The use of AIC for ARMA model order selection has theoretical support, particularly for forecasting, as ARMA models inspired the original AIC paper Akaike (1974). However, without proper likelihood maximization, a strategy that considers only a single parameter initialization may not truly minimize AIC. In this framework, likelihood maximization and over-parameterization are interconnected: all candidate models must be maximized for likelihood, or users risk selecting over-parameterized models that fail to minimize the intended information criterion. For the current study, the choice of AIC versus other popular information criteria such as the corrected AIC (AICC) or Bayesian information criterion (BIC) is unimportant: all of these approaches rely on proper optimization of the likelihood function, which is the problem we are addressing here.

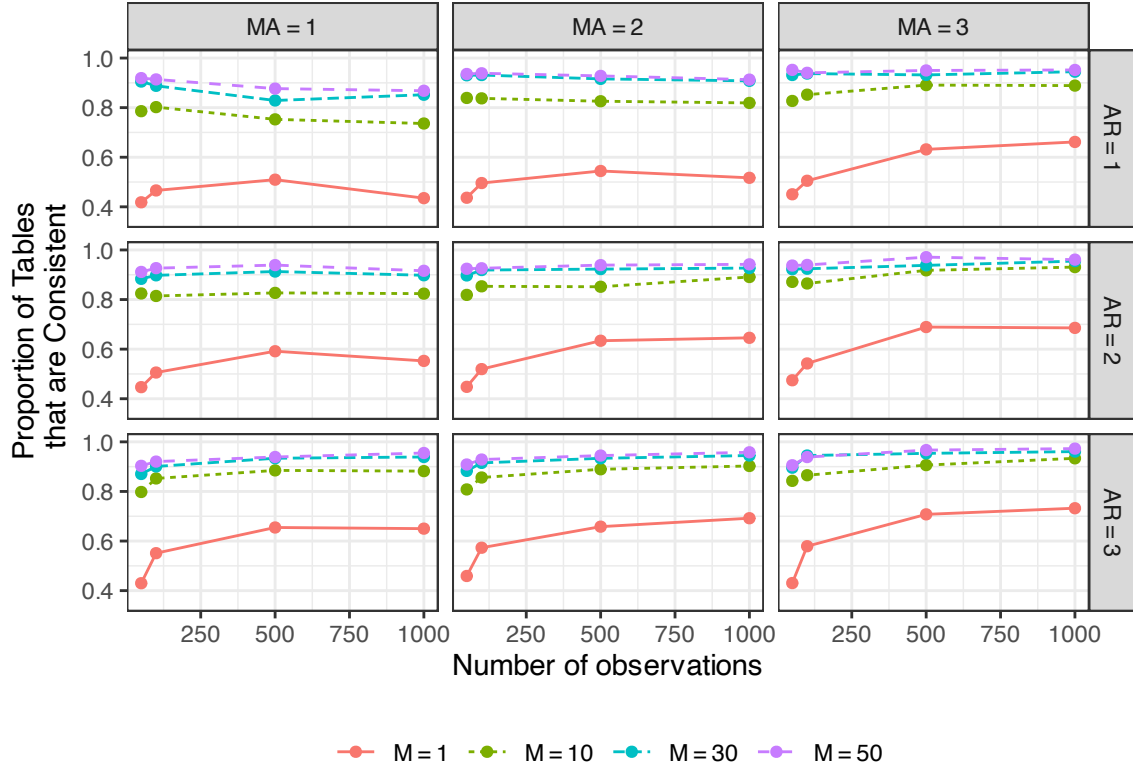


Figure 2.4: Data is generated from $\text{ARMA}(p, q)$ models with $(p, q) \in \{1, 2, 3\}^2$, and the corresponding AIC table is created. The Y-axis shows the percentage of tables that were consistent. M is the number of times a maxima is observed before the algorithm terminates, so $M = 1$ corresponds to the standard maximization procedure.

Classical ARMA modeling addresses this by recommending diagnostic plots to determine appropriate model order and advising against simultaneously adding AR and MA components. In this approach, the additional difficulty in parameter estimation associated with fitting models containing parameter redundancies is avoided by not fitting overly complex models when possible. Despite this, shortcomings in likelihood maximization can occur even in models without parameter redundancies (Figs 2.1, 2.2, and 2.4), necessitating the exploration of multiple parameter initializations. Further, the increasing preference of using automated software to pick the model size using an information criterion suggests the importance of using software that reliably maximizes model likelihoods even in the presence of over-parameterization.

2.3 Annual Depths of Lake Michigan

In this example, we illustrate how improperly maximized likelihoods can lead to inconsistencies and uncertainty in a real data analysis scenario. Additionally, we show how the common practice of using the estimated standard error for calibrated parameters can misleadingly support the inclusion of model parameters. We consider a dataset containing annual observations on the average depth of Lake Michigan-Huron, recorded the first day of each year from 1860-2014 (Fig 2.5) NOAA (2016). We wish to develop an ARMA model for these data, which is a standard task in time series analysis Shumway and Stoffer (2017).

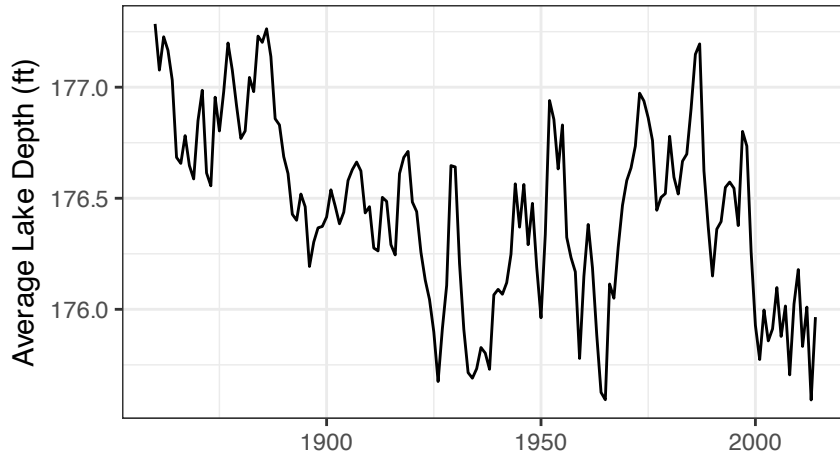


Figure 2.5: Average depth of Lake Michigan-Huron from 1860-2014.

Diagnostic tests, such as sample autocorrelation and normal quantile plots for residuals, suggest that it is reasonable to model the data in Fig 2.5 as a weakly stationary Gaussian $\text{ARMA}(p, q)$ process for some non-negative integers p and q . While an ARIMA model may also be reasonable for these data, we first consider fitting an ARMA model because we would like to avoid the possibility of over-differencing the data. The next step is to determine appropriate values of p and q ; after some initial investigation, multiple combinations of p and q seem plausible, and therefore we decide to choose the values of p and q that minimize the AIC. For simplicity, we create a table of AIC values for all possible combinations of $(p, q) \in \{0, 1, 2, 3\}^2$ (Table 2.1). Using the AIC as the model selection criterion, the selected model size is $\text{ARMA}(2, 1)$.

Recall that the AIC is defined as:

$$\text{AIC} = -2 \max_{\psi} \ell(\psi; x^*) + 2d, \quad (2.8)$$

where $\ell(\psi; x^*)$ denotes the log-likelihood of a model indexed by parameter vector $\psi \in \mathbb{R}^d$,

(a) Single parameter initialization.

	MA0	MA1	MA2	MA3
AR0	166.8	46.6	7.3	-15.0
AR1	-38.0	-37.4	-35.5	-33.8
AR2	-37.3	-38.4	-36.9	-34.9
AR3	-35.5	-35.2	-33.0	-33.3

(b) Multiple parameter initializations.

	MA0	MA1	MA2	MA3
AR0	166.8	46.6	7.3	-15.0
AR1	-38.0	-37.4	-35.5	-33.8
AR2	-37.3	-38.4	-36.9	-34.9
AR3	-35.5	-36.9	-36.4	-36.2

Table 2.1: AIC values for an ARMA(p, q) model fit to Lake Michigan-Huron depths. Table 2.1a was computed using only a single parameter initialization. Table 2.1b was computed using Algorithm 1. Highlighted cells show where the likelihood was improved (AIC reduced) using our algorithm.

$d \geq 1$, given the observed data x^* . In the case of an ARMA model with an intercept, $d = p + q + 2$, where the additional parameter corresponds to a variance estimate. If either p or q increases by one, then a corresponding increase in AIC values greater than two suggests that the *inclusion* of an additional parameter resulted in a *decrease* in the maximum of the log-likelihood, which is mathematically impossible under proper optimization. Several such cases are present in Table 2.1a, for example increasing from an ARMA(2, 2) model to a ARMA(3, 2) model results in a decrease of 1.0 log-likelihood units. In this case, using our multiple restart algorithm eliminates all instances of mathematical inconsistencies (Table 2.1b). We refer to tables that have log-likelihood values larger for any smaller nested model within the table as *inconsistent*.

Suppose a scientist is confronted with a mathematically implausible table of nominally maximized likelihoods (Table 2.1a). How much should they worry about this? Is it acceptable to publish scientific results that demonstrate a nominally maximized likelihood is not, in fact, maximized? Can researchers confidently trust the scientific implications of a fitted model if there is evidence of improper optimization in some of the candidate models? Given a choice, a researcher should prefer to use maximization algorithms reliable enough to make such situations rare. In the Lake Michigan example, improved estimation does not change which model is selected or the final parameter estimates, but it does remove inconsistencies that could lead to these concerns (Table 2.1b).

Minimizing the AIC (or an alternative information criterion) is not the only accepted approach to order selection. A classical perspective on model selection involves consulting sample autocorrelation plots, partial autocorrelation plots, conducting tests such as Ljung-Box over various lags, studying the polynomial roots of fitted models, and checking properties of the residuals of the fitted models Box and Jenkins (1970); Brockwell and Davis (1991); Shumway and Stoffer (2017). This approach helps avoid fitting models that are possibly

Table 2.2: Parameter values of an ARMA model fit to Lake Michigan-Huron data.

	ϕ_1	ϕ_2	θ_1	Intercept
Estimate	-0.053	0.791	1.000	176.460
s.e.	0.052	0.053	0.024	0.121

over-parameterized. However, additional computational power and increasing volumes of data have favored automated data analysis strategies that fit many models and evaluate them using a model selection criterion. In principle, a simple model selection criterion such as AIC can address parsimony and guard against over-parameterization as well. Diagnostic inspection can be combined with these automated approaches. For example, a table of AIC values can be generated, and models with promising likelihoods can be explored further Brockwell and Davis (1991).

When possible, there may be general agreement that the best approach is to combine modern computational resources with careful attention to model diagnostics, considering the data and the scientific task at hand. Improved maximization facilitates this process by eliminating distractions resulting from incomplete maximization.

2.3.1 Parameter uncertainty

Default output from fitting an ARMA model in R or Python includes estimates for parameter values and their standard errors, calculated using Fisher’s information matrix. If the ARMA model with the lowest AIC value is chosen to describe the Lake Michigan data, then an ARMA(2, 1) model is selected. The estimated coefficients and standard errors obtained after fitting this model are reported in Table 2.2. The small standard error for $\hat{\theta}_1$ reported in this table suggests a high-level of confidence that the parameter has a value near 1. Taken at face value, these estimates seem to strongly favor the inclusion of the MA(1) term in the model.

However, our simulation studies have suggested that these confidence intervals can be misleading, and that PLCIs are more reliable alternatives. The 95% PLCI for the parameter (Fig 2.6A) is much larger than the confidence interval created using these standard errors. The steep curve in the immediate vicinity of $\hat{\theta}_1$ may explain the small standard error estimates for this parameter and the corresponding tight confidence intervals created using Fisher’s identity matrix. Alternative evidence indicates the potential for nearly canceling roots (Fig 2.7), in which case the MA(1) term may not be needed in the model.

Both types of confidence intervals considered in this example rely on asymptotic justifications, but we can further investigate the finite sample properties using a simulation study. We fit both ARMA(2, 1) and AR(1) models to the data, and conduct a boot-strap simulation

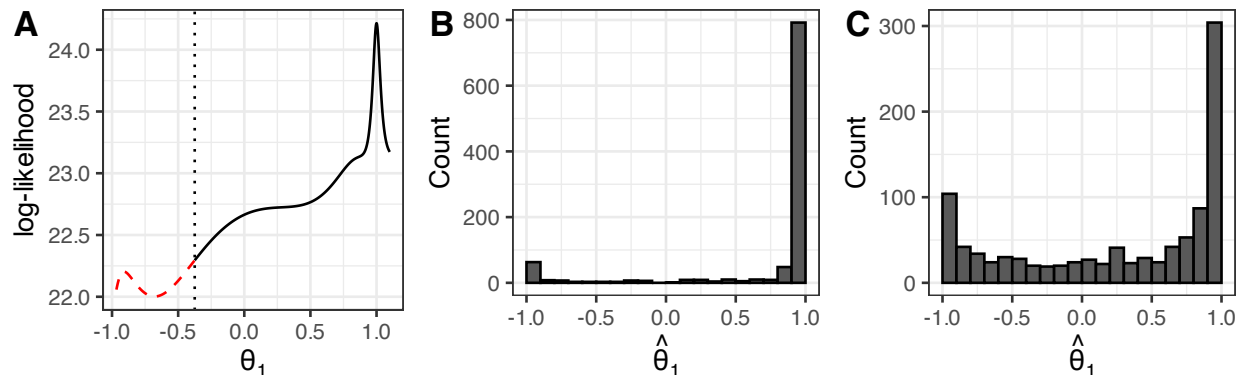


Figure 2.6: Evidence for an AR(1) model for the Lake Michigan-Huron data. (A) Profile likelihood confidence interval (PLCI) for θ_1 which includes the value $\theta_1 = 0$. The vertical dotted line represents the lower end of the approximate confidence interval; all points on the solid black line lie within the confidence interval, and points on the dashed red line are outside the interval. (B) Histogram of re-estimated θ_1 values using simulated data simulated from the ARMA(2,1) model that was calibrated to the Lake Michigan-Huron data. (C) Histogram of re-estimated θ_1 values using data simulated from the AR(1) model that was calibrated to the Lake Michigan-Huron data.

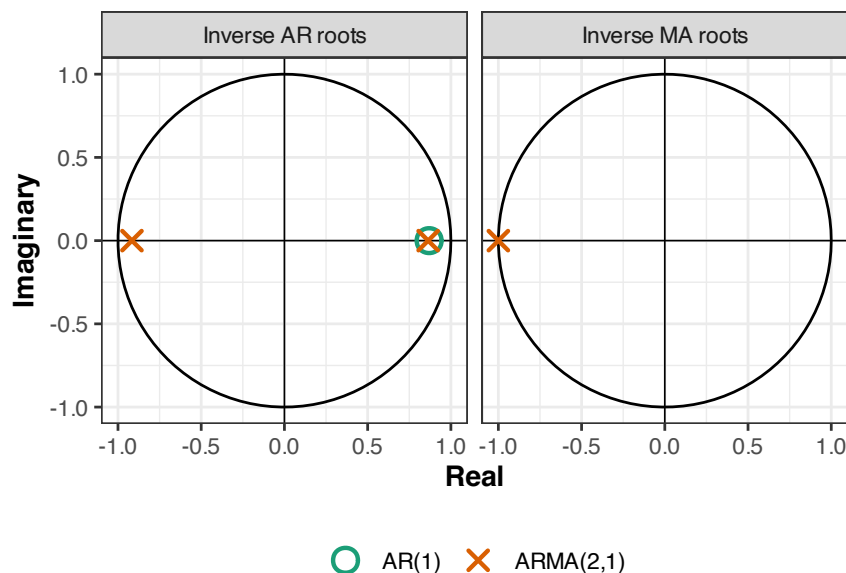


Figure 2.7: Inverted AR and MA polynomial roots to the fitted AR(1) and ARMA(2,1) models to the Lake Michigan-Huron data using a single parameter initialization.

study by simulating 1000 datasets from each of the fitted models. We then re-estimate an ARMA(2,1) model to each of these datasets and record the estimated coefficients. A histogram containing the estimated values of $\hat{\theta}_1$ when the data are generated from the

ARMA(2, 1) and AR(1) models fit to the Lake Huron-Michigan data are shown in Fig 2.6B and Fig 2.6C, respectively.

The shape of this histogram in Fig 2.6B mimics that of the profile log-likelihood surface in Fig 2.6A, confirming that a large confidence interval is needed in order to obtain a 95% confidence interval. In Fig 2.6C, a large number of $\hat{\theta}_1$ coefficients are estimated near 1 when the generating model is AR(1). Combining this result with the nearly canceling roots of the ARMA(2, 1) model (Fig 2.7), we cannot reject the hypothesis that the data were generated from a AR(1) model, even though the Fisher information standard errors suggest that the data should be modeled with a nonzero θ_1 coefficient.

2.4 Discussion

A significant motivation for our work is the observation that commonly used statistical software that purports to maximize ARMA model likelihoods fails to do so for a large number of examples. In addition to improving parameter estimates, proper maximization of ARMA model likelihoods is crucial because ARMA models are often used to model serial correlations in regression analyses. In this context, researchers may perform likelihood ratio hypothesis tests for regression coefficients, and the validity of these tests depends on proper likelihood optimization.

An important consequence of improved likelihood maximization is better model selection. A common approach to selecting an ARMA model involves fitting different sizes of models and choosing the one that minimizes an information criterion, such as the AIC. Fitting multiple models results in having a higher probability that at least one candidate model was not properly maximized. Since AIC assumes the parameters correspond to maximized likelihoods, enhancements in likelihood maximization can lead to different model selections. Consequently, methodology relying on existing estimation methods—like the popular `auto.arima` function in the `forecast` package in R Hyndman and Khandakar (2008), which minimizes the AIC of a group of candidate models without explicitly displaying an AIC table—will be impacted by improved estimates.

Our proposed algorithm is supported by existing theory on likelihood evaluation of linear state-space models via the Kalman Filter Kalman (1960), the same as the current existing standard approach for parameter estimation. The simulation studies that we have conducted, however, demonstrate the importance of considering multiple parameter initializations in order to fully maximize model likelihoods. These simulations provide a conservative estimate of how frequently our algorithm results in improved likelihoods compared to existing standards. A common situation where our algorithm is expected to provide even larger improvements

than those reported here is in the presence of missing data, a primary motivator of the likelihood maximization procedure of existing software Ripley (2002). In this situation, the well-informed CSS initialization is not available, and the default approach is to initialize at the origin, resulting in a greater need to attempt multiple parameter initializations.

Parameter estimates corresponding to higher likelihood values are not necessarily scientifically preferable to alternative regions of parameter space with lower likelihood values Le Cam (1990). Sometimes, our improved estimates may result in models with nearly canceling roots, parameters near boundary conditions, or otherwise unfavorable statistical properties. On other occasions, our method can rescue a naive optimization attempt from a local maximum having those unfavorable properties. Practitioners should carefully evaluate fitted models to ensure they are appropriate for the data and problem at hand.

The primary limitation of our approach is that it achieves higher likelihoods at the cost of processing speed, which is more pronounced with large datasets. However, our algorithm is most necessary for small datasets ($n \ll 10000$), where default parameter initialization strategies may perform poorly. Therefore, our algorithm is most beneficial for small to moderate sample sizes, where the additional computational cost is generally negligible. The compute time of our algorithm is approximately K times slower than the default approach, where K is the number of unique parameter initializations. This is only an approximation of the actual additional cost as not all initializations require the same amount of processing time in order to converge. In particular, initializations that are already close to local maximum will generally converge much quicker than those that are further away.

Our proposed algorithm for ARMA parameter estimation significantly advances statistical practice by addressing a frequently occurring optimization deficiency. Because existing software can also be leveraged to mitigate the issue, the largest contribution of this work may be highlighting the prevalence of this optimization problem. Traditional random initialization approaches software fail to uniformly cover the entire range of possible models and often produce many initializations outside the accepted range. Our algorithm offers a computationally efficient and practically convenient solution, providing a robust approach to parameter initialization and estimation that ensures adequate coverage of all possible models. We have shown that it provides a new standard for best practice in the field of time series analysis.

CHAPTER 3

Informing policy via dynamic models: Cholera in Haiti

3.1 Introduction

Regulation of biological populations is a fundamental topic in epidemiology, ecology, fisheries and agriculture. Population dynamics may be nonlinear and stochastic, with the resulting complexities compounded by incomplete understanding of the underlying biological mechanisms and by partial observability of the system variables. Quantitative models for these dynamic systems offer potential for designing effective control measures (Vandermeer and Goldberg, 2013; He et al., 2010). Developing and testing models for dynamic systems, and assessing their fitness for guiding policy, is a challenging statistical task (King et al., 2016). Questions of interest include: What indications should we look for in the data to assess whether the model-based inferences are trustworthy? What diagnostic tests and model variations can and should be considered in the course of the data analysis? What are the possible trade-offs of increasing model complexity, such as the inclusion of interactions across spatial units?

This case study investigates the use of dynamic models and spatiotemporal data to inform public health policy in the context of the cholera outbreak in Haiti, which started in 2010. Various dynamic models were developed to study this outbreak: searching PubMed with keywords “Haiti, cholera, model” we obtained 22 studies that utilized deterministic mechanistic dynamic models (Lee et al., 2020; Tuite et al., 2011; Andrews and Basu, 2011; Botelho et al., 2021; Fitzgibbon et al., 2020; Eisenberg et al., 2013; Rinaldo et al., 2012; Chao et al., 2011; Date et al., 2011; Lin et al., 2019; Abrams et al., 2013; Collins and Duffy, 2021; Akman and Schaefer, 2015; Trevisin et al., 2022; Mavian et al., 2020; Collins and Govinder, 2014; Kelly Jr et al., 2016; Capone et al., 2015; Leung et al., 2022; Mari et al., 2015; Gatto et al., 2012; Kühn et al., 2014) and 11 studies that used stochastic models (Kirpich et al., 2015; Lee et al., 2020; Pasetto et al., 2018; Mukandavire et al., 2013; Kirpich et al., 2017;

Lewnard et al., 2016; Kunkel et al., 2017; Mukandavire and Morris Jr, 2015; Sallah et al., 2017; Azman et al., 2012, 2015). Incidence data on the outbreak are available at various spatial scales, motivating 17 studies in our literature review to consider spatially explicit dynamic models (Lee et al., 2020; Tuite et al., 2011; Pasetto et al., 2018; Fitzgibbon et al., 2020; Eisenberg et al., 2013; Rinaldo et al., 2012; Chao et al., 2011; Abrams et al., 2013; Trevisin et al., 2022; Sallah et al., 2017; Collins and Govinder, 2014; Kelly Jr et al., 2016; Azman et al., 2012; Leung et al., 2022; Kühn et al., 2014; Mari et al., 2015; Gatto et al., 2012). Here we focus on a multi-group modeling exercise by Lee et al. (2020) in which four expert modeling teams developed models to the same dataset with the goal of comparing conclusions on the feasibility of eliminating cholera by a vaccination campaign. Model 1 is stochastic and describes cholera at the national level; Model 2 is deterministic with spatial structure, and includes transmission via contaminated water; Model 3 is stochastic with spatial structure, and accounts for measured rainfall. Model 4 has an agent-based construction, featuring considerable mechanistic detail but limited ability to calibrate these details to data. These modeling strategies were selected by Lee et al. (2020) to represent the range of approaches used in the research community. We focus on Models 1–3, as the strengths and weaknesses of the agent-based modeling approach (Tracy et al., 2018) are outside the scope of this article. In addition, agent-based models were less widely used, the agent based model in Lee et al. (2020) being the only model of this class that was found in our literature review. The data that were used in Lee et al. (2020), and that we reanalyze, are displayed in Fig. 3.1.

The four independent teams were given the task of estimating the potential effect of prospective oral cholera vaccine (OCV) programs. While OCV is accepted as a safe and effective tool for controlling the spread of cholera, the global stockpile of OCV doses remains limited (Pezzoli, 2020). Advances in OCV technology and vaccine availability, however, raised the possibility of planning a national vaccination program. The possibility of controlling the Haiti cholera outbreak via OCV was considered by various research groups (Lee et al., 2020; Ivers, 2017; Andrews and Basu, 2011; Walton et al., 2011; Matias et al., 2017; Chao et al., 2011; Fung et al., 2013; Date et al., 2011; Leung et al., 2022; Azman et al., 2015). In the Lee et al. (2020) study, certain data were shared between the groups, including demography and vaccination history; vaccine efficacy was also fixed at a shared value between groups. Beyond this, the groups made autonomous decisions on what to include and exclude from their models. Despite their autonomy, the four independent teams obtained a consensus that an extensive nationwide vaccination campaign would be necessary to eliminate cholera from Haiti, estimating that a large number of cumulative cholera cases would be observed in the absence of additional vaccination efforts (Figure 3 and 4 of Lee et al. (2020)). These forecasts are inconsistent with the prolonged period with no confirmed cholera cases between February,

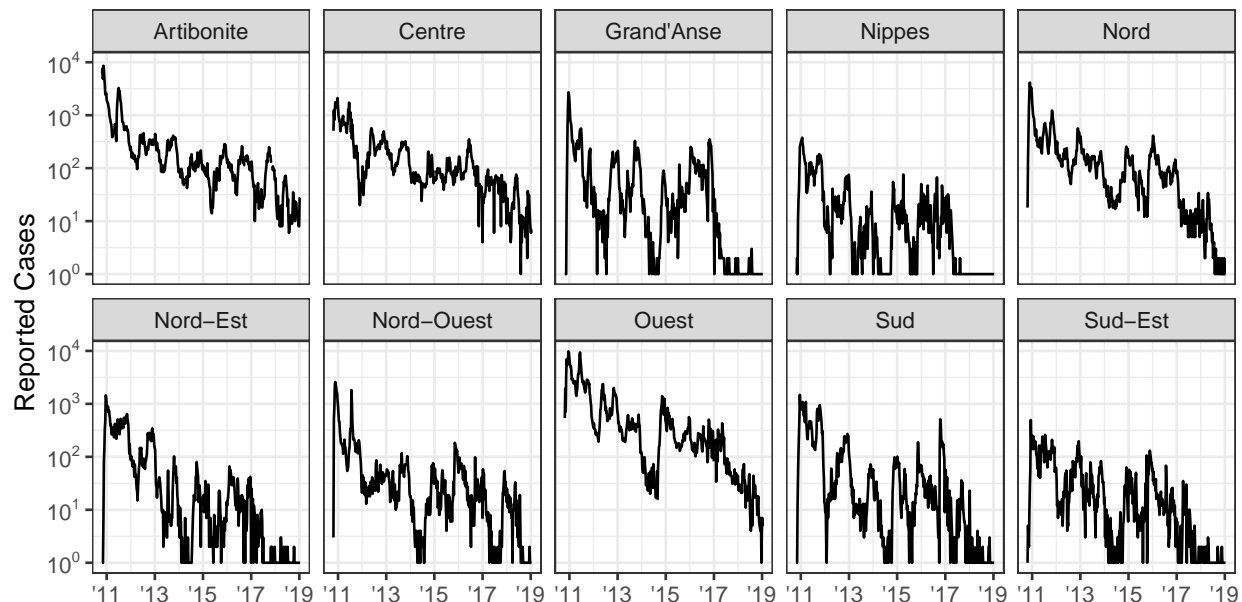


Figure 3.1: Weekly reported cholera cases in Haiti from October 2010 to January 2019 for each of the 10 administrative departments.

2019 and September, 2022 (Trevisin et al., 2022). Though cholera has recently reemerged in Haiti (Rubin et al., 2022; Pan American Health Organization, 2023), the inability to accurately forecast cholera incidence from 2019-2022 prompts us to consider retrospectively what may have been done differently in order to obtain more reliable conclusions, leading to recommendations for future studies.

The discrepancy between the model-based conclusions of Lee et al. (2020) and the prolonged absence of cholera in Haiti has been debated (Francois, 2020; Rebaudet et al., 2020; Henrys et al., 2020; Lee et al., 2020). Suggested origins of this discrepancy include the use of unrealistic models (Rebaudet et al., 2020) and unrealistic criteria for cholera elimination (Henrys et al., 2020). We find a more nuanced conclusion: attention to methodological details in model fitting, diagnosis and forecasting can improve each of the proposed model’s ability to quantitatively describe observed data. This improved ability may result in more accurate forecasts and facilitates the exploration of model assumptions. Based on this retrospective analysis, we offer suggestions on fitting mechanistic models to dynamic systems for future studies.

Numerous guidelines have been proposed for using mechanistic models to inform policy, reviewed by Behrend et al. (2020), who identify the importance of stakeholder engagement, transparency, reproducibility, uncertainty communication, and testable model outcomes.

These and related principles are echoed by other influential articles (Saltelli et al., 2020; Donnelly et al., 2018). Additional literature emphasizes model calibration and evaluation techniques (Dahabreh et al., 2017; Egger et al., 2017; Peñaloza Ramos et al., 2015). These guidelines often lack implementation specifics. As an example, Lee et al. (2020) largely adhere to the principles of Behrend et al. (2020)—though assessing the extent of stakeholder engagement is challenging—yet their projections are inconsistent with actual cholera incidence data from 2019 to 2022, demonstrating the limitations of current standards. We provide methodology for rigorous statistical calibration and evaluation of dynamic models (as advocated by Saltelli (2019)), thereby expanding on the prevailing guidance. We specifically emphasize principles that prove essential in our case study. Complementary methodological suggestions arising from a spatio-temporal analysis of COVID-19 are detailed in Li et al. (2023).

Our recommendations are presented in the context of a case study, with the goal of demonstrating how careful adherence to statistical principles may result in improved model fits. We proceed by introducing the general modeling scheme employed by Models 1–3 and provide details of each individual model; we then describe how each model is calibrated to data, and present a systematic approach to examining and refining these models. Specifically, we focus on how to develop and test variations of the proposed models, as well as diagnosing the models once they have been assimilated to incidence reports. This includes a comprehensive tutorial on performing inference with Model 3 (Appendix B.6), a highly non-linear, spatially explicit stochastic model, a challenging task that is possible due to recent methodological advancements. We then use the improved model fits to project cholera incidence in Haiti under various vaccination scenarios considered by Lee et al. (2020). Finally, we conclude with a discussion of the results, in which we relate our general recommendations for model based inference of biological systems to the case study of the Haiti cholera outbreak.

3.2 Materials and methods

3.2.1 Mechanistic models for disease modeling

Mechanistic models representing biological phenomena are valuable for epidemiology and consequently for public health policy (Lofgren et al., 2014; McCabe and Donnelly, 2021). More broadly, they have useful roles throughout biology, especially when combined with statistical methods that properly account for stochasticity and nonlinearity (May, 2004). In some situations, modern machine learning methods can outperform mechanistic models on epidemiological forecasting tasks (Lau et al., 2022; Baker et al., 2018). The predictive skill

of non-mechanistic models can reveal limitations in mechanistic models, but cannot readily replace the scientific understanding obtained by describing the biological dynamics of the system in a mathematical model (Baker et al., 2018; Prosperi et al., 2020).

In this article, we refer to models that focus on learning relationships between variables in a dataset as *associative*, whereas models that incorporate a known scientific property of the system we call *causal* or *mechanistic*. The danger in using forecasting techniques which rely on associative models to predict the consequence of interventions is called the Lucas critique in an econometric context. Lucas et al. (1976) pointed out that it is naive to predict the effects of an intervention on a given system based entirely on historical associations. To successfully predict the effect of an intervention, a model should therefore both provide a quantitative explanation of existing data and should have a causal interpretation: a manipulation of the system should correspond quantitatively with the corresponding change to the model. This motivates the development of mechanistic models, which provides a statistical fit to the available data while also supporting a causal interpretation. Despite their limited ability to project the effect of interventions on a system, associative models can be effectively used to make inference on an certain features of a system. In our literature review, there were 14 studies that used associative models to describe various aspects of the cholera epidemic (Hulland et al., 2019; Moise et al., 2020; Piarroux et al., 2011; Matias et al., 2017; Eisenberg et al., 2013; Rebaudet et al., 2019; Ivers et al., 2015; Raila and Anderson, 2017; Cuneo et al., 2017; Charles et al., 2014; Richterman et al., 2019,?; Bengtsson et al., 2015; Li et al., 2016).

The four mechanistic models of Lee et al. (2020) were deliberately developed with limited coordination. This allows us to treat the models as fairly independently developed expert approaches to understanding cholera transmission. However, it led to differences in notation, and in subsets of the data chosen for analysis, that hinder direct comparison. Here, we have created a common notational framework that facilitates model comparison, and put all comparable model parameters—including parameters that were estimated or held constant—into Table 3.1. Translations back to the original notation of Lee et al. (2020) are given in Table 3.2.

Each model describes the cholera dynamics as a partially observed Markov process (POMP) with a latent state vector $\mathbf{X}(t)$ for each continuous time point t . N observations on the system are collected at time points t_1, \dots, t_N , written as $t_{1:N}$. The observation at time t_n is modeled by the random vector \mathbf{Y}_n . While the latent process exists between observation times, the value of the latent state at observations times is of particular interest. We therefore write $\mathbf{X}_n = \mathbf{X}(t_n)$ to denote the value of the latent process at the n th observation time, and $\mathbf{X}_{1:N}$ is the collection of latent state values for all observed time points. The observable random variables $\mathbf{Y}_{1:N}$ are assumed to be conditionally independent given $\mathbf{X}_{0:N}$. Together,

Mechanism	Model 1	Model 2	Model 3
Infection (day)	$\mu_{IR}^{-1} = 2.0^{\dagger}$ (3.8)	$\mu_{IR}^{-1} = 7.0^{\dagger}$ (3.17)	$\mu_{IR}^{-1} = 5.0^{\dagger}$ (3.30)
Latency (day)	$\mu_{EI}^{-1} = 1.4^{\dagger}$ (3.7)	$\mu_{EI}^{-1} = 1.3^{\dagger}$ (3.16)	
Seasonality	$\beta_{1:6} = (1.4, 1.2, 1.1, 1.1, 1.4, 1.0)$ (3.4) $\zeta = -0.04^*$ (3.42)	$a = 0.4^{\dagger}$ (3.14) $\phi = 0.97^*$ (3.14)	$a = 1.00$ (3.34) $r = 0.78$ (3.34)
Immunity (yr)	$\mu_{RS}^{-1} = 8.0^{\dagger}$ (3.9)	$\mu_{RS}^{-1} = 1.4 \times 10^{11}$ (3.18) $\omega_1^{-1} = 1.0^{\dagger}$ (3.20) $\omega_2^{-1} = 5.0^{\dagger}$ (3.21)	$\mu_{RS}^{-1} = 8.0^{\dagger}$ (3.32)
Birth/death (yr ⁻¹)	$\mu_S = 10^{-2} \times 2.23^{\dagger}$ (3.11) $\delta = 10^{-3} \times 7.5^{\dagger}$ (3.11)		$\delta = 10^{-2} \times 1.59^{\dagger}$ (3.31) $\delta_C = 1.46^{\dagger}$ (3.34)
Sympt. frac.	$f_z(t) = c\vartheta^*(t - \tau_d)^{\dagger}$ (3.6-3.7)	$f = 0.2^{\dagger}$ (3.16)	$f = 0.25^{\dagger}$ (3.29)
Asympt. infectivity	$\epsilon = 0.05^{\dagger}$ (3.3)	$\epsilon = 0.001^{\dagger}$ (3.14) $\epsilon_W = 10^{-7}$ (3.22)	$\epsilon = 1^{\dagger}$ (3.27) $\epsilon_W = 0.008$ (3.34)
Human to human	$\beta_{1:6}$ as above (3.3)	$\beta = 5.97 \times 10^{-15} \text{ yr}^{-1}$ (3.14)	$\beta_{1:10} = (0.8, 0.0, 0.4, 0.2, 0.5, 0.5, 0.4, 0.1, 0.3, 0.1) \times 10^{-6} \text{ yr}^{-1}$ (3.27)
Water to human		$W_{\text{sat}} = 10^{5\dagger}$ (3.14) $\beta_W = 1.1 \text{ yr}^{-1}$ (3.14)	$\beta_{W_{1:10}} = (4.7, 21.0, 25.0, 27.1, 5.3, 30.7, 10.2, 1.0, 11.9, 12.8) \text{ yr}^{-1}$ (3.27)
Human to water		$\mu_W = 179 \text{ wk}^{-1}$ (3.22)	$\mu_W = 9.77 \times 10^{-7} \frac{\text{km}^2}{\text{wk}}$ (3.34)
Water survival (wk)		$\delta_W^{-1} = 3^{\dagger}$ (3.23)	$\delta_W^{-1} = 0.11$ (3.35)
Mixing exponent	$\nu = 0.98$ (3.3)		
Process noise (wk ^{1/2})	$\sigma_{\text{proc}} = (0.09, 0.12)^*$ (3.3)		$\sigma_{\text{proc}} = 0.218$ (3.29)
Reporting rate	$\rho = 0.679$ (3.13)	$\rho = 0.20^{\dagger}$ (3.26)	$\rho = 0.98$ (3.26)
Observation variance	$\psi = (279.1, 78.3)$ (3.13)	$\psi = 1.3$ (3.26)	$\psi = 88.6$ (3.41)
Initial Values	$I_{0,0} = 7298$ $E_{0,0} = 350$		$I_{0,0}^{3,4} = (21, 6)^*$
Hurricane Parameters			$\beta_{W_{3,9}}^{hm} = (36.88, 31.64)^*$ (3.27) $h_{3,9}^{hm} = (98.98, 58.43)^*$ (3.27)

Table 3.1: References to the relevant equation are given in parentheses. Parameters that were fixed and not calibrated using the data are indicated with † ; all fixed parameters values were chosen to match the fixed parameter values of Lee et al. (2020). Parameters that were added during our re-analysis and were not considered by Lee et al. are indicated with $*$. Confidence intervals for model parameters are given in Appendix B.4. Translations back into the notation of Lee et al. (2020) are given in Table 3.2.

Parameter	Our Notation	Lee et al. (2020a)		
		1	2	3
Reporting Rate	ρ	ρ	ρ	ϵ_1, ϵ_2
Mixing Coefficient	ν	ν	—	—
Measurement Over-Dispersion	ψ	τ	—	p
Birth Rate	μ_S	μ	—	—
Natural Mortality Rate	δ	δ	—	μ
Cholera Mortality Rate	δ_C	—	—	α
Latent Period	$1/\mu_{EI}$	$1/\sigma$	$1/\gamma_E$	—
Recovery Rate	μ_{IR}	γ	γ	γ
Loss of Immunity	μ_{RS}	α	σ	ρ
Symptomatic Ratio	f	$1 - \theta_0$	k	σ
Asymptomatic Relative Infectiousness	ϵ	$1 - \kappa$	red_β	—
Human-to-Water Shedding	μ_W	—	μ	θ_I
Asymptomatic Relative Shedding	ϵ_W	—	red_μ	θ_A/θ_I
Seasonal Amplitude	a	—	α_s	λ
Transmission	β	β	β	c
Water-to-Human	β_W	—	β_W	β
Bacteria Mortality	δ_W	—	δ	μ_β
Vaccination Efficacy	ϑ	θ_{vk}	$\theta_1, \theta_2, \theta_{15}, \theta_{25}$	η_{1d}, η_{2d}
Process Over-dispersion	σ_{proc}	—	—	σ_w

Table 3.2: Translations between our common notation and notation used by Lee et al (2020).

with the density for the initial value of the latent state $\mathbf{X}_0 = \mathbf{X}(t_0)$, each model defines a joint density $f_{\mathbf{X}_{0:N}, \mathbf{Y}_{1:N}}(\mathbf{x}_{0:N}, \mathbf{y}_{1:N}; \theta)$, where θ is a parameter vector that indexes the model. The observed data $\mathbf{y}_{1:N}^*$, along with the unobserved true value of the latent state, are modeled as a realization of this joint distribution.

Because of the probabilistic nature of both the unobserved latent state and the observable random variables, it is possible to consider various marginal and conditional densities of these two jointly random vectors. An important example is the marginal density of the observed random vector $\mathbf{Y}_{1:N}$, evaluated at the observed data $\mathbf{y}_{1:N}^*$, as shown in Eq. (3.1):

$$f_{\mathbf{Y}_{1:N}}(\mathbf{y}_{1:N}^*; \theta) = \int f_{\mathbf{X}_{0:N}, \mathbf{Y}_{1:N}}(\mathbf{x}_{0:N}, \mathbf{y}_{1:N}^*; \theta) d\mathbf{x}_{0:N}. \quad (3.1)$$

When treated as a function of the parameter vector θ , this marginal density is called the *likelihood function*, which is the basis of likelihood based statistical inference.

Using the conditional independence of $\mathbf{Y}_{1:N}$ given $\mathbf{X}_{0:N}$ and the Markov property of $\mathbf{X}_{0:N}$, the joint density can be re-factored into the useful form given in Eq. (3.2):

$$f_{\mathbf{X}_{0:N}, \mathbf{Y}_{1:N}}(\mathbf{x}_{0:N}, \mathbf{y}_{1:N}; \theta) = f_{\mathbf{X}_0}(\mathbf{x}_0; \theta) \prod_{n=1}^N f_{\mathbf{X}_n | \mathbf{X}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}; \theta) f_{\mathbf{Y}_n | \mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n). \quad (3.2)$$

This factorization is useful because it demonstrates that POMP models may be completely described using three parts: the *initialization model* for the latent states $f_{\mathbf{X}_0}(\mathbf{x}_0; \theta)$; the *one-step transition density*, or the *process model* $f_{\mathbf{X}_n | \mathbf{X}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}; \theta)$; and the *measurement model* $f_{\mathbf{Y}_n | \mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n)$. In the following subsections, we describe Models 1–3 in terms of these three components. The latent state vector $\mathbf{X}(t)$ for each model consists of individuals labeled as susceptible (S), infected (I), asymptotically infected (A), vaccinated (V), and recovered (R), with various sub-divisions sometimes considered. The observable random vector $\mathbf{Y}_{1:N}$ represents the random vector of cholera incidence data for each model; Models 2 and 3 have metapopulation structure, meaning that each individual is a member of a spatial unit, denoted by a subscript $u \in 1 : U$, in which case we denote the observed data for each unit using $\mathbf{y}_{1:N}^* = \mathbf{y}_{1:N,1:U}^*$. Here, the spatial units are the $U = 10$ Haitian administrative départements (henceforth anglicized as departments).

While the complete model description is scientifically critical, as well as necessary for transparency and reproducibility, the model details are not essential to our methodological discussions of how to diagnose and address model misspecification with the purpose of informing policy. A first-time reader may choose to skim through the rest of this section,

and return later. Additional details about the numeric implementation of these models are provided in Appendix B.2. While each of the dynamic models considered in this manuscript can be fully described using the mathematical equations provided in the following section, diagrams of dynamic systems can be helpful to understand the equations. For this reason, we provide flow chart diagrams for Models 1–3 in Figures B.1–B.3.

Model 1

The latent state vector $\mathbf{X}(t) = (S_z(t), E_z(t), I_z(t), A_z(t), R_z(t), z \in 0 : Z)$ describes susceptible, latent (exposed), infected (and symptomatic), asymptomatic, and recovered individuals in vaccine cohort z at time t . Here, $z = 0$ corresponds to unvaccinated individuals, and $z \in 1 : Z$ describes hypothetical vaccination programs. Each program z indexes differences in both the number of doses administered (one versus two doses per individual) and the round of vaccine administration, separating individuals into compartments with distinct dynamics based on vaccination status. The force of infection is

$$\lambda(t) = \left(\sum_{z=0}^Z I_z(t) + \epsilon \sum_{z=0}^Z A_z(t) \right)^\nu \frac{d\Gamma(t)}{dt} \beta(t) / N, \quad (3.3)$$

where $\beta(t)$ is a periodic cubic spline representation of seasonality, given in terms of a B-spline basis $\{s_j(t), j \in 1:6\}$ and parameters $\beta_{1:6}$ as

$$\beta(t) = \bar{\beta} \exp \left(\sum_{j=1}^6 \beta_j s_j(t) \right), \quad (3.4)$$

where $\bar{\beta} = 1 \text{ (wk)}^{-1}$ is a dimensionality constant. The process noise $d\Gamma(t)/dt$ is multiplicative Gamma-distributed white noise, with infinitesimal variance parameter σ_{proc}^2 . Lee et al. (2020) included process noise in Model 3 but not in Model 1, i.e., they fixed $\sigma_{\text{proc}}^2 = 0$. Gamma white noise in the transmission rate gives rise to an over-dispersed latent Markov process (Bretó and Ionides, 2011) which has been found to improve the statistical fit of disease transmission models (Stocks et al., 2020; He et al., 2010).

For any time point in $t_{1:N}$, the process model $f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}; \theta)$ is defined by describing how individuals move from one latent state compartment to another. Per-capita transition

rates are given in Eqs. (3.5)-(3.12):

$$\mu_{S_z E_z} = \lambda(t), \quad (3.5)$$

$$\mu_{E_z I_z} = \mu_{EI} (1 - f_z(t)), \quad (3.6)$$

$$\mu_{E_z A_z} = \mu_{EI} f_z(t), \quad (3.7)$$

$$\mu_{I_z R_z} = \mu_{A_z R_z} = \mu_{IR}, \quad (3.8)$$

$$\mu_{R_z S_z} = \mu_{RS}, \quad (3.9)$$

$$\mu_{S_0 S_z} = \mu_{E_0 E_z} = \mu_{I_0 I_z} = \mu_{A_0 A_z} = \mu_{R_0 R_z} = \eta_z(t), \quad (3.10)$$

$$\mu_{S_z \bullet} = \mu_{E_z \bullet} = \mu_{I_z \bullet} = \mu_{A_z \bullet} = \mu_{R_z \bullet} = \delta, \quad (3.11)$$

$$\mu_{\bullet S_0} = \mu_S, \quad (3.12)$$

where $z \in 0 : Z$. Here, μ_{AB} is a transition rate from compartment A to B . We have an additional demographic source and sink compartment \bullet modeling entry into the study population due to birth or immigration, and exit from the study population due to death or immigration. Thus, $\mu_{A\bullet}$ is a rate of exiting the study population from compartment A and $\mu_{\bullet B}$ is a rate of entering the study population into compartment B .

In Model 1, the advantage afforded to vaccinated individuals is an increased probability that an infection is asymptomatic. Conditional on infection status, vaccinated individuals are also less infectious than their non-vaccinated counterparts by a rate of $\epsilon = 0.05$ in Eq. (3.3). In Eqs. (3.7) and (3.6) the asymptomatic ratio for non-vaccinated individuals is set $f_0(t) = 0$, so that the asymptomatic route is reserved for vaccinated individuals. For $z \in 1 : Z$, the vaccination cohort z is assigned a time τ_z , and we take $f_z(t) = c \vartheta^*(t - \tau_z)$ where $\vartheta^*(t)$ is efficacy at time t since vaccination for adults, a step-function represented in Table S4 of Lee et al. (2020), and $c = (1 - (1 - 0.4688) \times 0.11)$ is a correction to allow for reduced efficacy in the 11% of the population aged under 5 years. Single and double vaccine doses were modeled by changing the waning of protection; protection was modeled as equal between single and double dose until 52 weeks after vaccination, at which point the single dose becomes ineffective.

The latent state vector $\mathbf{X}(t)$ is initialized by setting the counts for each compartment and vaccination scenario $z \neq 0$ as zero, and introducing initial-value parameters $I_{0,0}$ and $E_{0,0}$ such that $R_0(0) = 0$, $I_0(0) = \text{Pop} \times I_{0,0}$, $E_0(0) = \text{Pop} \times E_{0,0}$ and $S_0(0) = \text{Pop} \times (1 - I_{0,0} - E_{0,0})$, where Pop is the total population of Haiti. These parameters are estimated using the observed data, with parameter estimates provided in Table 3.1.

The measurement model describes reported cholera cases Y_n at the nationally aggregated level at each week n come from a negative binomial distribution, where only a fraction ρ of

new weekly cases are reported. Because vaccinated individuals are treated as asymptomatic, all reported cases are modeled as transitions from the exposed to the infected compartment, as described in Eq. (3.13).

$$Y_n \mid \Delta N_{E.I.}(n) = \delta \sim \text{NB}(\rho\delta, \psi), \quad (3.13)$$

where Y_n is the reported cholera cases at time $t_n \in t_1 : t_N$ and $\Delta N_{E.I.}(n)$ is the sum total of individuals across each vaccination compartment $z \in 1 : Z$ who moved from compartment E_z to I_z since observation t_{n-1} . Here, $\text{NB}(\rho\delta, \psi)$ denotes a negative binomial distribution with mean $\rho\delta$ and variance $\rho\delta\left(1 + \frac{\rho\delta}{\psi}\right)$.

Model 2

Susceptible individuals are in compartments $S_{uz}(t)$, where $u \in 1 : U$ corresponds to the $U = 10$ departments, and $z \in 0 : 4$ describes vaccination status:

$z = 0$: Unvaccinated or waned vaccination protection.

$z = 1$: One dose at age under five years.

$z = 2$: Two doses at age under five years.

$z = 3$: One dose at age over five years.

$z = 4$: Two doses at age over five years.

Like Model 1, the process model $f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}; \theta)$ is primarily defined via the description of movement of individuals between compartments, however Model 2 also includes a dynamic description of a latent bacterial compartment as well. Individuals can progress to a latent infection E_{uz} followed by symptomatic infection I_{uz} with recovery to R_{uz} or asymptomatic infection A_{uz} with recovery to R_{uz}^A . The force of infection depends on both direct transmission and an aquatic reservoir, $W_u(t)$, and is given by

$$\lambda_u(t) = 0.5\left(1 + a \cos(2\pi t + \phi)\right) \frac{\beta_W W_u(t)}{W_{\text{sat}} + W_u(t)} + \beta \left\{ \sum_{z=0}^4 I_{uz}(t) + \epsilon \sum_{z=0}^4 A_{uz}(t) \right\}. \quad (3.14)$$

The latent state is therefore described by the vector $\mathbf{X}(t) = (S_{uz}(t), E_{uz}(t), I_{uz}(t), A_{uz}(t), R_{uz}(t), R_{uz}^A(t), W_u, u \in 1 : U, z \in 0 : 4)$. The cosine term in Eq. (3.14) accounts for annual seasonality, with a phase parameter ϕ . The Lee et al. (2020) implementation of Model 2 fixes $\phi = 0$.

Individuals move from department u to v at rate T_{uv} , and aquatic cholera moves at rate T_{uv}^W . The nonzero transition rates are

$$\mu_{S_{uz}E_{uz}} = (1 - \vartheta_z) \lambda_u(t), \quad (3.15)$$

$$\mu_{E_{uz}I_{uz}} = f\mu_{EI}, \quad \mu_{E_{uz}A_{uz}} = (1 - f)\mu_{EI}, \quad (3.16)$$

$$\mu_{I_{uz}R_{uz}} = \mu_{A_{uz}R_{uz}^A} = \mu_{IR}, \quad (3.17)$$

$$\mu_{R_{uz}S_{uz}} = \mu_{R_{uz}^A S_{uz}} = \mu_{RS}, \quad (3.18)$$

$$\mu_{S_{uz}S_{vz}} = \mu_{E_{uz}E_{vz}} = \mu_{I_{uz}I_{vz}} = \mu_{A_{uz}A_{vz}} = \mu_{R_{uz}R_{vz}} = \mu_{R_{uz}^A R_{vz}^A} = T_{uv}, \quad (3.19)$$

$$\mu_{S_{u1}S_{u0}} = \mu_{S_{u3}S_{u0}} = \omega_1, \quad (3.20)$$

$$\mu_{S_{u2}S_{u0}} = \mu_{S_{u4}S_{u0}} = \omega_2, \quad (3.21)$$

$$\mu_{\bullet W_u} = \mu_W \left\{ \sum_{z=0}^4 I_{uz}(t) + \epsilon_W \sum_{z=0}^4 A_{uz}(t) \right\}, \quad (3.22)$$

$$\mu_{W_u \bullet} = \delta_W, \quad (3.23)$$

$$\mu_{W_u W_v} = w_r T_{uv}^W. \quad (3.24)$$

In Eq. (3.19) the spatial coupling is specified by a gravity model,

$$T_{uv} = v_{\text{rate}} \times \frac{\text{Pop}_u \text{Pop}_v}{D_{uv}^2}, \quad (3.25)$$

where Pop_u is the mean population for department u , D_{uv} is a distance measure estimating average road distance between randomly chosen members of each population, and $v_{\text{rate}} = 10^{-12} \text{ km}^2 \text{ yr}^{-1}$ was fixed at the value used in Lee et al. (2020). In Eq. (3.24), T_{uv}^W is a measure of river flow between departments. The unit of $W_u(t)$ is cells per ml, with dose response modeled via a saturation constant of W_{sat} in Eq. (3.14). In Eq. (3.15), ϑ_z denotes the vaccine efficacy for each vaccination campaign $z \in Z$, with $\vartheta_0 = 0$, $\vartheta_1 = 0.429q$, $\vartheta_2 = 0.519q$, $\vartheta_3 = 0.429$, and $\vartheta_4 = 0.519$. Here, $q = 0.4688$ represents the reduced efficacy of the vaccination for children under the age of five years, and the values 0.429 and 0.519 are the median effectiveness of one and two doses over their effective period respectively, according to Table S4 in the supplement material of Lee et al. (2020). Because vaccine efficacy remains constant, individuals in this model transition from a vaccinated compartment to the susceptible compartment at the end of the vaccine coverage period.

The starting value for each element of the latent state vector $\mathbf{X}(0)$ are set to zero except for $I_{u0}(0) = y_u^*(0)/\rho$ and $R_{u0}(0) = \text{Pop}_u - I_{u0}(0)$, where $y_u^*(0)$ is the reported number of cholera cases in department u at time $t = 0$. Reported cases are described using a log-normal distribution, with the log-scale mean equal to the reporting rate ρ times the number of newly

infected individuals, following Eq. (3.26).

$$\log(Y_{u,n} + 1) \mid \Delta N_{E_u \cdot I_u}(n) = \delta_u \sim N(\log(\rho\delta_u + 1), \psi^2), \quad (3.26)$$

where $\Delta N_{E_u \cdot I_u}(n)$ is the sum total of individuals across vaccination compartment $z \in 0 : 4$ within unit u who moved from compartment E_{uz} to I_{uz} since observation t_{n-1} . Therefore, because the natural logarithm of observed case counts (plus one, to avoid taking the logarithm of zero) has a normal distribution, $Y_{u,n} + 1$ is assumed to follow a log-normal distribution with log-mean parameter $\log(\rho\Delta N_{E_u \cdot I_u}(n) + 1)$ and log-variance ψ^2 . We note that fitting a model with this measurement model is equivalent to fitting using least squares, with $\log(Y_{u,n} + 1)$ as the response variable. This measurement model differs from that used by Lee et al. (2020), who fit the model in two stages: epidemic and endemic phases, accounting for symptomatic infections in the endemic phase and both symptomatic and asymptomatic infections in the endemic period.

Model 3

The latent state is described as $\mathbf{X}(t) = (S_{uz}(t), I_u(t), A_u(t), R_{uzk}(t), W_u(t), u \in 0:U, z \in 0:4, k \in 1:3)$. Here, $z = 0$ corresponds to unvaccinated, $z = 2j - 1$ corresponds to a single dose on the j th vaccination campaign in unit u and $z = 2j$ corresponds to receiving two doses on the j th vaccination campaign. $k \in 1:3$ models non-exponential duration in the recovered class before waning of immunity. The processes model $f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}; \theta)$ describes the movement of individuals between latent compartments, as well as the birth and death process of local, unobserved bacterial compartments $W_u(t)$. The force of infection is

$$\lambda_u(t) = (\beta_{W_u} + 1_{(t \geq t_{hm})} \beta_{W_u}^{hm} e^{-h_u^{hm}(t-t_{hm})}) \frac{W_u(t)}{1 + W_u(t)} + \beta_u \sum_{v \neq u} (I_v(t) + \epsilon A_v(t)), \quad (3.27)$$

where t_{hm} is the time Hurricane Matthew struck Haiti (Ferreira, 2016), and $1_{(A)}$ is the indicator function for event A . In Lee et al. (2020), $\beta_{W_u}^{hm}$ and h_u^{hm} were set to zero for all u ; the need to account for the effect Hurricane Matthew had on cholera transmission for this model is explored in Sec. S5 of the supplement.

Per-capita transition rates are used for both compartments representing human counts

and the aquatic reservoir of bacteria; these rates are given in Eqs. (3.28)–(3.35).

$$\mu_{S_{uz}I_u} = f \lambda_u (1 - \vartheta_{uz}(t)) d\Gamma/dt, \quad (3.28)$$

$$\mu_{S_{uz}A_u} = (1 - f) \lambda_u (1 - \vartheta_{uz}(t)) d\Gamma/dt, \quad (3.29)$$

$$\mu_{I_u R_{uz1}} = \mu_{A_u R_{uz1}} = \mu_{IR}, \quad (3.30)$$

$$\mu_{I_u S_{u0}} = \delta + \delta_C, \quad \mu_{A_u S_{u0}} = \delta \quad (3.31)$$

$$\mu_{R_{uz1} R_{uz2}} = \mu_{R_{uz2} R_{uz3}} = 3\mu_{RS}, \quad (3.32)$$

$$\mu_{R_{uzk} S_{u0}} = \delta + 3\mu_{RS} \mathbf{1}_{\{k=3\}}, \quad (3.33)$$

$$\mu_{\bullet W_u} = \left[1 + a(J_u(t))^r\right] \text{Den}_u \mu_W [I_u(t) + \epsilon_W A_u(t)], \quad (3.34)$$

$$\mu_{W_u \bullet} = \delta_W. \quad (3.35)$$

As with Model 1, $d\Gamma_u(t)/dt$ is multiplicative Gamma-distributed white noise in Eqs. (3.28) and (3.29). In Eq. (3.34), $J_u(t)$ is a dimensionless measurement of precipitation that has been standardized by dividing the observed rainfall at time t by the maximum recorded rainfall in department u during the epidemic, and Den_u is the population density. Demographic stochasticity is accounted for by modeling non-cholera related death rate δ in each compartment, along with an additional death rate δ_C in Eq. (3.31) to account for cholera induced deaths among infected individuals. All deaths are balanced by births into the susceptible compartment in Eqs. (3.31) and (3.33), thereby maintaining constant population in each department.

Similar to Model 1, there are no distinct compartments for individuals under five years of age, and the vaccination efficacy is taken as a age adjusted weighted average of the efficacy for individuals both over and under five years of age: $\vartheta_{uz}(t) = c\vartheta^*(t - \tau_{uz})$, where τ_{uz} is the vaccination time for unit u and vaccination campaign z . The value c and the function ϑ^* are equivalent to those described in the Model 1 description.

The latent states of this model are initialized by enforcing the model dynamics on the incidence data from the start of the recorded cases until time t_0 , requiring that some of the available data be used to determine the initial values of the latent states. This is the same approach that was taken by Lee et al. (2020), who used the value $t_0 = 2014-02-22$; this choice of t_0 results in modeling roughly only 60% of the available data, some of which is later discarded for alternative reasons. We do not see any immediate reason that this model could not be extended to cover a larger range of the data, and chose the value $t_0 = 2010-11-13$.

This choice of t_0 corresponds to using approximately 1% of the available data to determine initial values of the latent states. In addition to modeling a larger portion of the available data, this choice of t_0 corresponds to an important real-world event, as daily reports from

each of the departments were not being sent to the health ministry until November 10, 2010 (Barzilay et al., 2013); this choice of t_0 therefore makes $\mathbf{Y}(t_1)$ the first week of data once daily reports are being sent to the health ministry. The few observation times that exist before t_0 are used to initialize the model by enforcing model dynamics on these preliminary observations. For convenience, we denote these observations as t_{-3}, t_{-2} and t_{-1} ; as before, we let $y_{u,-k} = Y_u(t_{-k})$ denote the observed case count for unit u at time point t_{-k} , where $k \in 1 : 3$. Equations for the initial values of non-zero latent states are provided in Eqs. (3.36)–(3.40); these equations match those that were used by Lee et al. (2020), the primary change being a change to the value of t_0 .

$$I_u(t_0) = \frac{y_{u,-1}^*}{7\rho(\mu_{IR} + (\delta + \delta_C)/365)}, \quad (3.36)$$

$$A_u(t_0) = \frac{I_{u0}(t_0)(1-f)}{f}, \quad (3.37)$$

$$R_{u01}(t_0) = R_{u02}(t_0) = R_{u03}(t_0) = \left(\frac{\sum_{k=-3}^0 y_{u,k}^*}{\rho f} - (I_{u0}(t_0) + A_{u0}(t_0)) \right) / 3 \quad (3.38)$$

$$S_{u0}(t_0) = \text{Pop}_u - I_u(t_0) - A_u(t_0) - \sum_{k=1}^3 R_{u0k}(t_0) \quad (3.39)$$

$$W_u(t_0) = [1 + a\tilde{J}^r] \text{Den}_u \mu_W [I_u(t_0) + \epsilon_W A_{u0}(t)] / \mu_W. \quad (3.40)$$

In Eq. (3.40), $\tilde{J} = 0.002376$ is the median adjusted rainfall over the observation period. One important consideration to make with this parameter initialization model is when $y_{u,-1}^* = 0$, which occurs for units $u \in \{3, 4\}$, which correspond to the Grand'Anse and Nippes departments, respectively. When this is the case, each of the infectious $I_u(t_0)$, asymptomatic $A_u(t_0)$, and bacterial reservoir $W_u(t_0)$ compartments have a value of zero. Because Model 3 models cholera transmission primarily by means of the bacterial reservoir, this makes it nearly impossible for an outbreak to occur. Therefore for units $u \in \{3, 4\}$, we introduce initial value parameters $I_{0,0}^3$ and $I_{0,0}^4$, and calibrate these parameter values using the data. The resulting parameter estimates are used to obtain the remaining non-zero initial values of the latent states using Eqs. (3.37)–(3.40).

Finally, reported cholera cases are assumed to stem from individuals who develop symptoms and seek healthcare. Therefore reported cases are assumed to come from an over-dispersed negative binomial model, given the increase in infected individuals:

$$Y_{u,n} \mid \Delta N_{S_u \cdot I_u}(t) = \delta_u \sim \text{NB}(\rho\delta_u, \psi), \quad (3.41)$$

where $\Delta N_{S_u, I_{uz}}(n)$ is the number of individuals who moved from compartment S_{uz} to I_u since observation t_{n-1} . This measurement model is a minor change from that used by Lee et al. (2020), which allowed for a change in the reporting rate on January 1st, 2018. Their estimated reporting rate value before and after this period were 0.97 and 0.097, respectively. An instantaneous change from near perfect to almost non-existent reporting can be problematic, as it forces the model to explain the observed reduction in reported cases as a decrease in the reporting of cases, rather than a decrease in the prevalence of cholera. To avoid this extreme possibility, we do not allow a change in reporting rate when fitting the model.

3.2.2 Model Fitting

Each of the three models considered in this study describes cholera dynamics as a partially observed Markov process (POMP) (King et al., 2016), with the understanding that the deterministic Model 2 is a special case of a Markov processes solving a stochastic differential equation in the limit as the noise parameter goes to zero. Each model is indexed by a parameter vector, θ , and different values of θ can result in qualitative differences in the predicted behavior of the system. Therefore, the choice of θ used to make inference about the system can greatly affect model-based conclusions (Saltelli et al., 2020). Elements of θ can be fixed at a constant value based on scientific understanding of the system, but parameters can also be calibrated to data by maximizing a measure of congruency between the observed data and the model’s mechanistic structure. Calibrating model parameters to observed data does not guarantee that the resulting model successfully approximates real-world mechanisms, since the model description of the dynamic system may be incorrect and does not change as the model is calibrated to data. However, the congruency between the model and observed data serves as a proxy for the congruency between the model and the true underlying dynamic system. As such, it is desirable to obtain the best possible fit of the proposed mechanistic structure to the observed data.

In this article we follow Lee et al. (2020) by calibrating the parameters of each of our models using maximum likelihood, as described in Eq. (3.1). The likelihood for each of the fitted models—and the corresponding AIC values for model comparisons that include an adjustment for the number of calibrated parameters—is provided in Table 3.3. In the following subsections we describe in detail our approach to calibrating the three proposed mechanistic models to observed cholera incidence data. The main alternative to maximum likelihood estimation is Bayesian inference via Markov chain Monte Carlo, used to analyze the Haiti cholera epidemic by Andrews and Basu (2011); Pasetto et al. (2018); Rinaldo et al. (2012); Lewnard et al. (2016); Trevisin et al. (2022); Sallah et al. (2017); Azman et al. (2012);

Kühn et al. (2014); Mari et al. (2015); Gatto et al. (2012).

	Model 1	Model 2	Model 3
Log-likelihood	−2728.1 (−3030.9) ¹	−21957.3 (−29367.4)	−17332.9 (−33832.6) ²
Number of Fit Parameters	15 (20)	6 (6)	34 (29)
AIC	5486.3 (6101.8) ¹	43926.5 (58746.9)	34733.9 (67723.2) ²
Benchmark AIC	5585.3	36961.0	35945.2

Table 3.3: AIC values for Models 1–3 and their benchmarks. Values in parentheses are corresponding values obtained using the models of Lee et al. (2020). ¹The reported likelihood is an upper bound of the likelihood of the model in Lee et al. (2020). ²In Lee et al. (2020), Model 3 was fit to a subset of the data (March 2014 onward, excluding data from Ouest in 2015–2016). On this subset, their model has a likelihood of −9721.2. On this same subset, our model has a likelihood of −7219.5. Details of estimating the likelihood of the models used in Lee et al. (2020) are provided in Appendix B.5.

Calibrating Model 1 Parameters

Model 1 is a highly nonlinear over-dispersed stochastic dynamic model, favoring a scientifically plausible description of cholera dynamics rather than one that is statistically convenient (He et al., 2010). This results in the inability to obtain a closed form expression of the joint model density—described in Eq. (3.2). Therefore in order to perform likelihood based inference on this model, we are restricted to use parameter estimation techniques that have the *plug-and-play* property, which is that the fitting procedure only requires the ability to simulate the latent process rather than evaluating transition densities (Bretó et al., 2009; He et al., 2010); in the context of the notation and definitions employed in this article, this means that we only require the ability to simulate from $f_{\mathbf{X}_0}(\mathbf{x}_0; \theta)$ and $f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}; \theta)$ rather than needing to evaluate these densities. Plug-and-play algorithms include Bayesian approaches like ABC and PMCMC (Toni et al., 2009; Andrieu et al., 2010), but here we use algorithms that enable maximum likelihood estimation. To our knowledge, the only plug-and-play methods that have been effectively used to maximize the likelihood for arbitrary nonlinear POMP models are iterated filtering algorithms (Ionides et al., 2015), which modify the well-known *particle filter* (Arulampalam et al., 2002).

The particle filter, also referred to as sequential Monte Carlo, is a simulation based method that is frequently used in Bayesian inference to approximate the posterior distribution of

latent states. This algorithm can also be used to accurately approximate the log-likelihood of a POMP model, defined as the integral in Eq. (3.1). Iterated filtering algorithms, such as IF2 (Ionides et al., 2015), extend the particle filter by performing a random walk for each parameter and particle; these perturbations are carried out iteratively over multiple filtering operations, using the collection of parameters from the previous filtering pass as the parameter initialization for the next iteration, and decreasing the random walk variance at each step. With a sufficient number of iterations, the resulting parameter values converge to a region of the parameter space that maximizes the model likelihood.

The ability to maximize the likelihood allows for likelihood-based inference, such as performing statistical tests for potential model improvements. We demonstrate this capability by proposing a log-linear trend ζ in transmission in Eq. (3.4):

$$\beta(t) = \bar{\beta} \exp \left(\sum_{j=1}^6 \beta_s s_j(t) + \zeta \bar{t} \right), \quad (3.42)$$

where $\bar{t} = \frac{t - (t_N + t_0)/2}{t_N - (t_N + t_0)/2}$, so that $\bar{t} \in [-1, 1]$. The proposal of a trend in transmission is a result of observing an apparent decrease in reported cholera infections from 2012-2019 in Fig. 3.1. While several factors may contribute to this decrease, one explanation is that case-area targeted interventions (CATIs), which included education sessions, increased monitoring, household decontamination, soap distribution, and water chlorination in infected areas (Rebaudet et al., 2019), may have substantially reduced cholera transmission over time (Rebaudet et al., 2021).

We perform a statistical test to determine whether or not the data indicate the presence of a trend in transmissibility. To do this, we perform a profile-likelihood search on the parameter ζ and obtain a 95% confidence interval via a Monte Carlo Adjusted Profile (MCAP) (Ionides et al., 2017). Lee et al. (2020) implemented Model 1 by fitting two distinct phases: an epidemic phase from October 2010 through March 2015, and an endemic phase from March 2015 onward. We similarly allow the re-estimation of process and measurement overdispersion parameters (σ_{proc}^2 and ψ), and require that the latent Markov process $X(t)$ carry over from one phase into the next. The resulting 95% confidence interval for ζ is $(-0.098, -0.009)$, with the full results displayed in Fig. 3.2. These results are suggestive that the inclusion of a trend in the transmission rate improves the quantitative ability of Model 1 to describe the observed data. The maximum likelihood estimate for ζ corresponds to a 7.3% reduction to the transmission rate over the course of the outbreak, with a 95% confidence interval of (1.8%, 17.9%) for the overall reduction in transmission. The reported results for Model 1 in the remainder of this article were obtained with the inclusion of the parameter ζ . The inclusion of a trend in transmission rate demonstrates a class of model variation that can be

highly beneficial to consider: the model variation has a plausible scientific justification, and is easily testable using likelihood based methods.

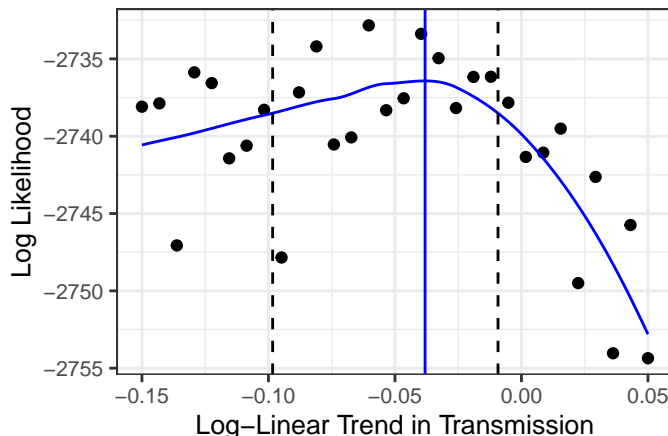


Figure 3.2: Monte Carlo adjusted profile (MCAP) of ζ for Model 1. The blue curve is the MCAP, the vertical blue line indicates the MLE, and the vertical dashed lines indicate the 95% confidence interval.

If a mechanistic model including a feature (such as a representation of a mechanism, or the inclusion of a covariate) fits better than mechanistic models without that feature, and also has competitive fit compared to associative benchmarks, this may be taken as evidence supporting the scientific relevance of the feature. As for any analysis of observational data, we must be alert to the possibility of confounding. For a covariate, this shows up in a similar way to regression analysis: the covariate under investigation could be a proxy for some other unmodeled phenomenon or unmeasured covariate.

The statistical evidence of a trend in transmission rate in this model could be explained by any trending variable (such as hygiene improvements, or changes in population behavior), resulting in confounding from collinear covariates. Alternatively, it is possible that the negative trend observed in the incidence data could be attributed to a decreasing reporting rate rather than decreasing transmission rate. This could be formally tested by comparing models with either trend specification. We did not do this because evidence suggests that reporting rate was maintained or increased (Figure 1 of Rebaudet et al. (2021)). We instead argue that a decreasing transmission rate is a plausible way to explain the decrease in cases over time, as there is alternative evidence that supports this model (Rebaudet et al., 2019, 2021; Michel et al., 2019). It is not practical to test all remotely plausible model variations, yet a strongly supported conclusion should avoid ruling out untested hypotheses. The robust statistical conclusion for our analysis is that a model which allows for change fits better than one which does not, and a trend in transmission is a plausible way to do this.

We implemented Model 1 using the **pomp** package (King et al., 2016), relying heavily on the source code provided by Lee et al. (2020). Both analyses used the **mif2** implementation of the IF2 algorithm to estimate θ by maximum likelihood. One change we made in the statistical analysis that led to larger model likelihoods was increasing the computational effort in the numerical maximization. While IF2 enables parameter estimation for a large class of models, the theoretic ability to maximize the likelihood depends on asymptotics in both the number of particles and the number of filtering iterations. Many Monte Carlo replications are then required to quantify and further reduce the error. The large increase in the log-likelihood for Model 1 (Table 3.3) can primarily be attributed to increasing the computational effort used to calibrate the model. This result highlights the importance of carefully determining the necessary computational effort needed to maximize model likelihoods and acting accordingly. In this case study, this was done by performing standard diagnostics for the IF2 and particle filter algorithms (King et al., 2016; Li et al., 2023; Pons-Salort and Grassly, 2018; Laneri et al., 2010). Given the considerable computational costs of simulation-based algorithms, we find it useful to perform an initial assessment using hyperparameter values—such as the number of particles, filtering iterations, and replicates based on different parameter initializations—that enable relatively quick calculations. The insights obtained from this preliminary analysis help in accurately determining the amount of computation that is required to achieve reliable outcomes. Simulations from the initial conditions of our fitted model are plotted against the observed incidence data in Fig. 3.3.

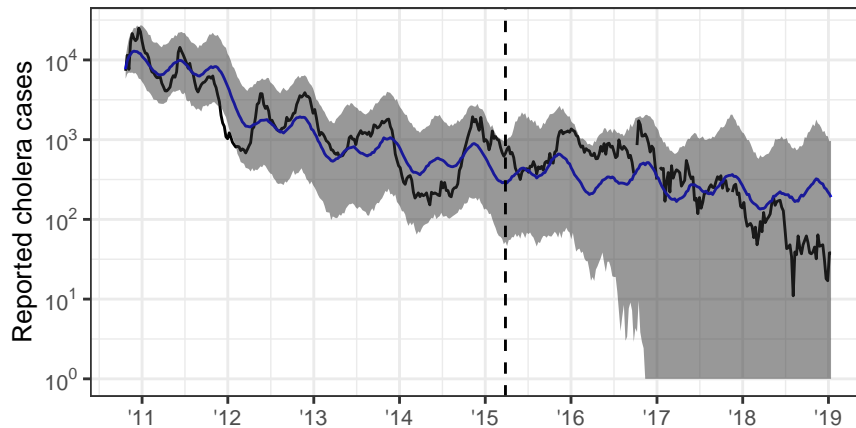


Figure 3.3: The black curve is observed data, the blue curve is median of 500 simulations from initial conditions using estimated parameters, and the vertical dashed line represents break-point when parameters are refit.

Calibrating Model 2 Parameters

Model 2 is a deterministic compartmental model defined by a set of coupled differential equations. The use of deterministic compartment models have a long history in the field of infectious disease epidemiology (Kermack and McKendrick, 1927; Brauer, 2017; Giordano et al., 2020), and can be justified by asymptotic considerations in a large-population limit (Dadlani et al., 2020; Ndii and Supriatna, 2017). Because the process model of Model 2 is deterministic, maximum likelihood estimation reduces to a least squares calculation when combined with a Gaussian measurement model (Eq. (3.26)). Lee et al. (2020) fit two versions of Model 2 based on a presupposed change in cholera transmission from an epidemic phase to endemic phase that occurred in March, 2014. The inclusion of a change-point in model states and parameters increased the flexibility of the model and hence the ability to fit the observed data. The increase in model flexibility, however, resulted in hidden states that were inconsistent between model phases. The inclusion of a model break-point by Lee et al. (2020) is perhaps due to a challenging feature of fitting a deterministic model via least squares: discrepancies between model trajectories and observed case counts in highly infectious periods of a disease outbreak will result in greater penalty than the discrepancies between model trajectories and observed case counts in times of relatively low infectiousness. This results in a bias towards accurately describing periods of high infectiousness. This bias is particularly troublesome for modeling cholera dynamics in Haiti: the inability to accurately fit times of low infectiousness may result in poor model forecasts, as few cases of cholera were observed in the last few years of the epidemic.

To combat this issue, we fit the model to log-transformed case counts, since the log scale stabilizes the variation during periods of high and low incidence. An alternative solution is to change the measurement model to include overdispersion, as was done in Models 1 and 3. This permits the consideration of demographic stochasticity, which is dominant for small infected populations, together with log scale stochasticity (also called multiplicative, or environmental, or extra-demographic) which is dominant at high population counts. Here we chose to fit the model to transformed case counts rather than adding overdispersion to the measurement model with the goal of minimizing the changes to the model proposed by Lee et al. (2020).

We implemented this model using the `spatPomp` R package (Asfaw et al., 2024). The model was then fit using the subplex algorithm (King and Rowan, 2020). A comparison of the trajectory of the fitted model to the data is given in Fig. 3.4.

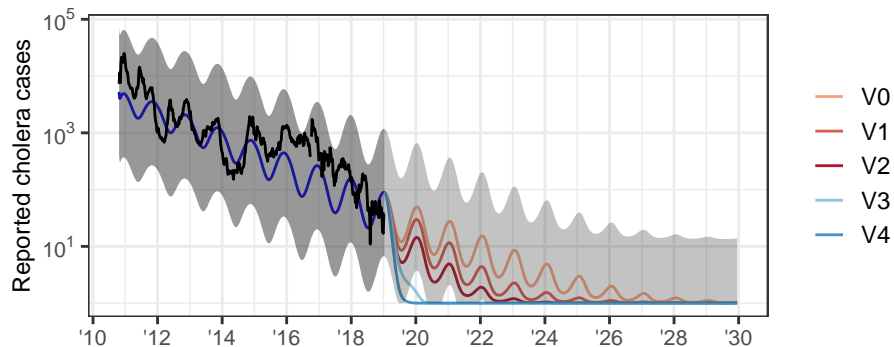


Figure 3.4: The black line shows the nationally aggregated weekly cholera incidence data. The blue curve from 2012-2019 is the trajectory of the calibrated version of Model 2. Projections under the various vaccination scenarios, which are discussed in detail in the **Forecasts** subsection are also included. The gray ribbons represent a 95% interval obtained from the log-normal measurement model. To avoid over-plotting, measurement variance is only plotted for the V0 vaccination scenario.

3.2.2.1 Calibrating Model 3 Parameters

Model 3 describes cholera dynamics in Haiti using a metapopulation model, where the hidden states in each administrative department has an effect on the dynamics in other departments. The decision to address metapopulation dynamics using a spatially explicit model, rather than to aggregate over space, is double-edged. Evidence for the former approach has been provided in previous studies (King et al., 2015), including the specific case of heterogeneity between Haitian departments in cholera transmission (Collins and Govinder, 2014). However, a legitimate preference for simplicity can support a decision to consider nationally aggregated models (Saltelli et al., 2020; Green and Armstrong, 2015).

In our literature review, 17 articles considered dynamic models that incorporate spatial heterogeneity (Lee et al., 2020; Tuite et al., 2011; Pasetto et al., 2018; Fitzgibbon et al., 2020; Eisenberg et al., 2013; Rinaldo et al., 2012; Chao et al., 2011; Abrams et al., 2013; Trevisin et al., 2022; Sallah et al., 2017; Collins and Govinder, 2014; Kelly Jr et al., 2016; Azman et al., 2012; Leung et al., 2022; Kühn et al., 2014; Mari et al., 2015; Gatto et al., 2012). All but four (Lee et al., 2020; Pasetto et al., 2018; Sallah et al., 2017; Azman et al., 2012) of these studies used deterministic dynamic models: this greatly simplifies the process of calibrating model parameters to incidence data, though deterministic models can struggle to describe complex stochastic dynamics. The model in (Pasetto et al., 2018) was fit using an Ensemble Kalman Filter (EnKF) (Evensen, 2009); though EnKF scales favorably with the number of spatial units, it relies on linearization of latent states which can be problematic for highly nonlinear systems (Evensen et al., 2022; Ionides et al., 2021). Alternative approaches used

to fit stochastic models included making additional simplifying assumptions to aid in the fitting process (Lee et al., 2020), and using MCMC algorithms (Sallah et al., 2017; Azman et al., 2012) which require specific structures in the latent dynamics, making these algorithms non plug-and-play. In this subsection, we present how the recently developed iterated block particle filter (IBPF) algorithm (Ning and Ionides, 2023; Ionides et al., 2022) can be used to fit a spatially explicit stochastic dynamic model to incidence data.

One issue that arises when fitting spatially explicit models is that parameter estimation techniques based on the particle filter become computationally intractable as the number of spatial units increases. This is a result of the approximation error of particle filters growing exponentially in the dimension of the model (Rebeschini and van Handel, 2015; Park and Ionides, 2020). To avoid the approximation error present in high-dimensional models, Lee et al. (2020) simplified the problem of estimating the parameters of Model 3 by creating an approximate version of the model where the units are independent given the observed data. Reducing a spatially coupled model to individual units in this fashion requires special treatment of any interactive mechanisms between spatial units, such as found in Eq. (3.27). Because the simplified, spatially-decoupled version of Model 3 implemented in Lee et al. (2020) relies on the observed cholera cases, the calibrated model cannot readily be used to obtain forecasts. Therefore, in order to obtain model forecasts, Lee et al. (2020) used the parameters estimates from the spatially-decoupled approximation of Model 3 to obtain forecasts using the fully coupled version of the model. This approach of model calibration and forecasting avoids the issue of particle depletion, but may also be problematic. One concern is that cholera dynamics in department u are highly related to the dynamics in the remaining departments; calibrating model parameters while conditioning on the observed cases in other departments may therefore lead to an over-dependence on observed cholera cases. Another concern is that the two versions of the model are not the same, resulting in sub-optimal parameter estimates for the spatially coupled model, as parameters that maximize the likelihood of the decoupled model almost certainly do not maximize the likelihood of the fully coupled model. These two concerns may explain the unrealistic forecasts and low likelihood of Model 3 in Lee et al. (2020) (Table 3.3).

At the time Lee et al. (2020) conducted their study, there was no known algorithm that could readily be used to maximize the likelihood of an arbitrary meta-population POMP model with coupled spatial dynamics, which justifies the spatial decoupling approximation that was used to calibrate model parameters. For our analysis, we calibrate the parameters of the spatially coupled version of Model 3 using the IBPF algorithm (Ionides et al., 2022). This algorithm extends the work of Ning and Ionides (Ning and Ionides, 2023), who provided theoretic justification for the version of the algorithm that only estimates unit-specific

parameters. The IBPF algorithm enables us to directly estimate the parameters of models describing high-dimensional partially-observed nonlinear dynamic systems via likelihood maximization. The ability to directly estimate parameters of Model 3 is responsible for the large increase in model likelihoods reported in Table 3.3. Simulations from the fitted model are displayed in Fig. 3.5.

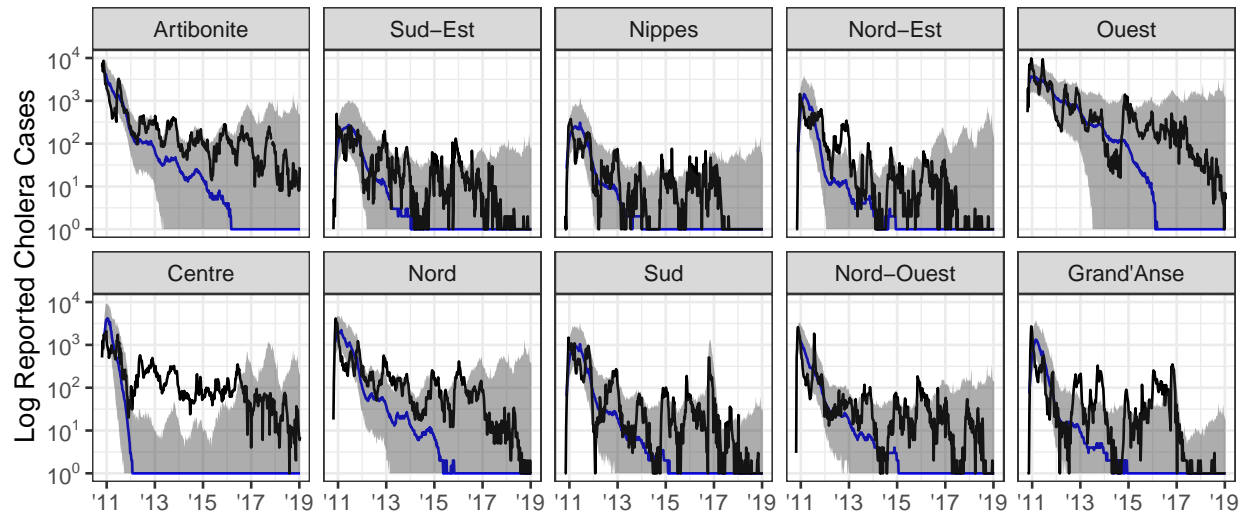


Figure 3.5: Simulations from initial conditions using the spatially coupled version of Model 3. The black curve represents true case count, the blue line the median of 500 simulations from the model, and the gray ribbons representing 95% confidence interval.

3.2.3 Model Diagnostics

The goal of parameter calibration—whether done using Bayesian or frequentist methods—is to find the best description of the observed data in the context of the model. Obtaining the best fitting set of parameters for a given model does not, however, guarantee that the model provides an accurate representation of the system under investigation. Model misspecification, which may be thought of as the omission of a mechanism in the model that is an important feature of the dynamic system, is inevitable at all levels of model complexity. To make progress, while accepting proper limitations, one must bear in mind the much-quoted observation of George Box (Box, 1979) that “all models are wrong but some are useful.” Beyond being good practical advice for applied statistics, this assertion is relevant for the philosophical justification of statistical inference as severe testing (Mayo, 2018). In this section, we discuss some tools for diagnosing mechanistic models with the goal of making the subjective assessment of model “usefulness” more objective. To do this, we will rely on

the quantitative ability of the model to match the observed data, which we call the model's *goodness-of-fit*, with the guiding principle that a model which cannot adequately describe observed data may not be reliable for useful purposes. Goodness-of-fit may provide evidence supporting the causal interpretation of one model versus another, but cannot by itself rule out the possibility of alternative explanations.

One common approach to assess a mechanistic model's goodness-of-fit is to compare simulations from the fitted model to the observed data. Visual inspection may indicate defects in the model, or may suggest that the observed data are a plausible realization of the fitted model. While visual comparisons can be informative, they provide only a weak and informal measure of the goodness-of-fit of a model. The study by Lee et al. (2020) provides an example of this: their models and parameter estimates resulted in simulations that visually resembled the observed data, yet resulted in model likelihoods that were considerably smaller than likelihoods that can be achieved (see Table 3.3). Alternative forms of model validation should therefore be used in conjunction with visual comparisons of simulations to observed data.

Another approach is to compare a quantitative measure of the model fit (such as MSE, predictive accuracy, or model likelihood) among all proposed models. These comparisons, which provide insight into how each model performs relative to the others, are quite common (Rinaldo et al., 2012; Sallah et al., 2017). To calibrate relative measures of fit, it is useful to compare against a model that has well-understood statistical ability to fit data, and we call this model a *benchmark*. Standard statistical models, interpreted as associative models without requiring any mechanistic interpretation of their parameters, provide suitable benchmarks. Examples include linear regression, auto regressive moving average (ARMA) time series models, or even independent and identically distributed measurements. Benchmarks enable us to evaluate the goodness of fit that can be expected of a suitable mechanistic model.

Associative models are not constrained to have a causal interpretation, and typically are designed with the sole goal of providing a statistical fit to data. Therefore, we should not require a candidate mechanistic model to beat all benchmarks. However, a mechanistic model which falls far short against benchmarks is evidently failing to explain some substantial aspect of the data. A convenient measure of fit should have interpretable differences that help to operationalize the meaning of far short. Ideally, the measure should also have favorable theoretical properties. Consequently, we focus on log-likelihood as a measure of goodness of fit, and we adjust for the degrees of freedom of the models to be compared by using the Akaike information criterion (AIC) (Akaike, 1974).

In some cases, a possible benchmark model could be a generally accepted mechanistic model, but often no such model is available. Because of this, we use a simple negative

binomial model with an auto regressive mean as our associative benchmark; this model is described in (3.43).

$$Y_n|Y_{n-1} \sim \text{NB}(\alpha + \beta Y_{n-1}, \varphi), \quad (3.43)$$

where $E(Y_n|Y_{n-1}) = \alpha + \beta Y_{n-1}$, and $\text{Var}(Y_n|Y_{n-1}) = E(Y_n|Y_{n-1}) + E(Y_n|Y_{n-1})^2 / \varphi$. To obtain a benchmark for models with a meta-population structure, we fit independent auto-regressive negative binomial models to each spatial unit. Under the assumption of independence, the log-likelihood of the benchmark on the entire collection of data can be obtained by summing up the log-likelihood for each independent model. In general, a spatially explicit model may not have well-defined individual log-likelihoods, and, in this case, comparisons to benchmarks must be made at the level of the joint model.

In the case where the case counts are large, an alternative benchmark recommended by He et al. (2010) is a log-linear Gaussian ARMA model; the theory and practice of ARMA models is well developed, and these linear models are appropriate on a log scale due to the exponential growth and decay characteristic of biological dynamics. We use the auto regressive negative binomial model, however, because the large number of weeks with zero recorded cholera cases in department level data makes a benchmark based on a continuous distribution problematic. Log-likelihoods and AIC values of Models 1–3 and of their respective benchmark models are provided in Table 3.3. Models that are fit to the same datasets can be directly compared using AIC values, making it a useful tool to compare to benchmark models. Though Models 2 and 3 are both fit to department level incidence reports, their AIC values are not directly comparable due to the way Model 3 initializes latent states.

It should be universal practice to present measures of goodness of fit for published models, and mechanistic models should be compared against benchmarks. In our literature review of the Haiti cholera epidemic, no non-mechanistic benchmark models were considered in any of the 32 papers that used dynamic models to describe cholera in order to obtain scientific conclusions. Including benchmarks would help authors and readers to detect and confront any major statistical limitations of the proposed mechanistic models. In addition, the published goodness of fit provides a concrete point of comparison for subsequent scientific investigations. When combined with online availability of data and code, objective measures of fit provide a powerful tool to accelerate scientific progress, following the paradigm of the *common task framework* (Donoho, 2017, Sec. 6).

The use of benchmarks may also be beneficial when developing models at differing spatial scales, where a direct comparison between model likelihoods is meaningless. In such a case, a benchmark model can be fit to each spatial resolution being considered, and each model

compared to their respective benchmark. Large advantages (or shortcomings) in model likelihood relative to the benchmark for a given spatial scale that are not present in other spatial scales may provide weak evidence for (or against) the statistical fit of models across a range of spatial resolutions.

Comparing model log-likelihoods to a suitable benchmark may not be sufficient to identify all the strengths and weaknesses of a given model. Additional techniques include the inspection of conditional log-likelihoods of each observation given the previous observations in order to understand how well the model describes each data point (Appendix B.6). Other tools include plotting the effective sample size of each observation (Liu, 2001); plotting the values of the hidden states from simulations (Appendix B.6); and comparing summary statistics of the observed data to simulations from the model (Wood, 2010; King et al., 2015).

3.2.4 Corroborating Fitted Models with Scientific Knowledge

The resulting mechanisms in a fitted model can be compared to current scientific knowledge about a system. Agreement between model-based inference and our current understanding of a system may be taken as a confirmation of both model-based conclusions and our scientific understanding. On the other hand, comparisons may generate unexpected results that have the potential to spark new scientific knowledge (Ganusov, 2016).

In the context of our case study, we demonstrate how the fit of Model 1 corroborates other evidence concerning the role of rainfall in cholera epidemics. Specifically, we examine the results of fitting the flexible cubic spline term in Model 1 (Eqs. (3.3)–(3.4)). The cubic splines permit flexible estimation of seasonality in the force of infection, $\beta(t)$. Fig. 3.6 shows that the estimated seasonal transmission rate β mimics the rainfall dynamics in Haiti, despite Model 1 not having access to rainfall data. This is consistent with previous studies that incorporated rainfall as an important part of their mechanistic model or otherwise argue that rainfall is an important driver of cholera dynamics in Haiti (Hulland et al., 2019; Kirpich et al., 2017; Lee et al., 2020; Moise et al., 2020; Pasetto et al., 2018; Kirpich et al., 2015; Eisenberg et al., 2013; Rinaldo et al., 2012; Mavian et al., 2020). The estimated seasonality also features an increased transmission rate during the fall, which was noticed at an earlier stage of the epidemic (Rinaldo et al., 2012). The high transmission rate in the fall may be a result of the increase transmission that occurred in the fall of 2016, when hurricane Matthew struck Haiti (Ferreira, 2016).

For any model-based inference, it is important to recognize and assess the modeling simplifications and assumptions that were used in order to arrive at the conclusions. In epidemiological studies, for example, quantitative understanding of individual-level processes

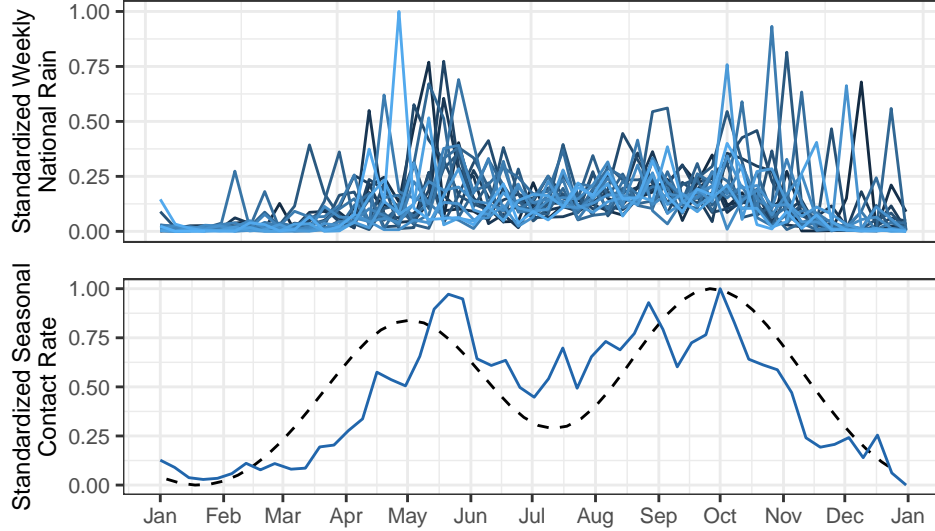


Figure 3.6: (Top) weekly rainfall in Haiti, lighter colors representing more recent years. (Bottom) estimated seasonality in the transmission rate (dashed line) plotted alongside mean rainfall (solid line). The outsized effect of rainfall in the fall may be due to Hurricane Matthew, which struck Haiti in October of 2016 and resulted in an increase of cholera cases in the nation.

may not perfectly match model parameters that were fit to population-level case counts, even when the model provides a strong statistical fit (He et al., 2010). This makes direct interpretation of estimated parameters delicate.

Our case study provides an example of this in the parameter estimate for the duration of natural immunity due to cholera infection, μ_{RS}^{-1} . Under the framework of Model 2, the best estimate for this parameter is 1.4×10^{11} yr, suggesting that individuals have permanent immunity to cholera once infected. Rather than interpreting this as scientific evidence that individuals have permanent immunity from cholera, this result suggests that Model 2 favors a regime where reinfection events are a negligible part of the dynamics. The depletion of susceptible individuals may be attributed to confounding mechanisms—such as localized vaccination programs and non-pharmaceutical interventions that reduce cholera transmission (Trevisin et al., 2022; Rebaudet et al., 2021)—that were not accounted for in the model. Perhaps the best interpretation of the estimated parameter, then, is that under the modeling framework that was used, the model most adequately describes the observed data by having a steady decrease in the number of susceptible individuals. The weak statistical fit of Model 2 compared to a log-linear benchmark (see Table 3.3) cautions us against drawing quantitative conclusions from this model. A model that has a poor statistical fit may nevertheless provide a useful conceptual framework for thinking about the system under investigation. However, a

claim that the model has been validated against data should be reserved for situations where the model provides a statistical fit that is competitive against alternative explanations.

A model which aspires to provide quantitative guidance for assessing interventions should provide a quantitative statistical fit for available data. However, strong statistical fit does not guarantee a correct causal structure: it does not even necessarily require the model to assert a causal explanation. A causal interpretation is strengthened by corroborative evidence. For example, reconstructed latent variables (such as numbers of susceptible and recovered individuals) should make sense in the context of alternative measurements of these variables (Grad et al., 2012). Similarly, parameters that have been calibrated to data should make sense in the context of alternative lines of evidence about the phenomena being modeled, while making allowance for the possibility that the interpretations of parameters may vary when modeling across differing spatial scales.

In Appendix B.6, we explore in more detail the process of model fitting and diagnostics for Model 3. Here we demonstrate that the model outperforms its benchmark model on the aggregate scale. However, when focusing on the spatial units with the highest incidence of cholera, Model 3 performs roughly the same as a simple benchmark. By comparing simulations from the fitted model to the filtering distribution, we see that the reconstructed latent states of the model favor higher levels of cholera transmission than what is typically observed in the incidence data. These results hint at the possibility of model misspecification, and warrant a degree of caution in interpreting the model’s outputs.

3.3 Results

3.3.1 Forecasts

Forecasts are an attempt to provide an accurate estimate of the future state of a system based on currently available data, together with an assessment of uncertainty. Forecasts from mechanistic models that are compatible with current scientific understanding may also provide estimates of the future effects of potential interventions. Further, they may enable real-time testing of new scientific hypotheses (Lewis et al., 2022).

Forecasts of a dynamic system should be consistent with the available data. It is particularly important that forecasts are consistent with the most recent information available, as recent data is likely to be more relevant than older data. While this assertion may seem self-evident, it is not the case for deterministic models, for which the initial conditions together with the parameters are sufficient for forecasting, and so recent data may not be consistent with model trajectories. Epidemiological forecasts based on deterministic models are not

uncommon in practice, despite their limitations (King et al., 2015). Lee et al. (2020) chose to obtain forecasts from all of their models by simulating forward from initial conditions, rather than conditioning forecasts based on the available data. This decision is possibly as a result of using a deterministic model, as forecasts from different models may only be considered comparable if they are obtained in the same way, which is most easily done by simulating from initial conditions because Model 2 is deterministic.

In contrast, for non-deterministic Models 1 and 3, we obtain forecasts by simulating future values using latent states that are harmonious with the most recent data. This is done by simulating forward from latent states drawn at the last observation time (t_N) from the filtering distribution $f_{\mathbf{x}_N|\mathbf{y}_{1:N}}(\mathbf{x}_N|\mathbf{y}_{1:N}^*; \hat{\theta})$. The decision to obtain model forecasts from initial conditions partially explains the unsuccessful forecasts of Lee et al. (2020). Table S7 in their supplement material, which contains results that were not discussed in their main article, shows that the subset of their simulations with zero cholera cases from 2019-2020 also correspond with its disappearance until 2022. These results support our argument that forecasts should be made by ensuring the starting point for the forecast is consistent with available data.

Uncertainty in just a single parameter can lead to drastically different forecasts (Saltelli et al., 2020). Therefore, parameter uncertainty should also be considered when obtaining model forecasts to influence policy. If a Bayesian technique is used for parameter estimation, a natural way to account for parameter uncertainty is to obtain simulations from the model where each simulation is obtained using parameters drawn from the estimated posterior distribution. For frequentist inference, one possible approach is obtaining model forecasts from various values of θ , where the values of θ are sampled proportionally according to their corresponding likelihoods (King et al., 2015). This approach is described in more detail in Appendix B.7. Both of these approaches share the similarity that parameters are chosen for the forecast approximately in proportion to their corresponding value of the likelihood function, $f_{\mathbf{y}_{1:N}}(\mathbf{y}_{1:N}^*; \theta)$. In this analysis, we do not construct forecasts accounting for parameter uncertainty as our focus is on the estimation and diagnosis of mechanistic models, rather than providing forecasts intended to influence policy. Furthermore, we use the projections from a single point estimate to highlight the deficiency of deterministic models that the only variability in model projections is a result of parameter and measurement uncertainty, which can lead to over-confidence in forecasts (King et al., 2015).

The primary forecasting goal of Lee et al. (2020) was to investigate the potential consequences of vaccination interventions on a system to inform policy. One outcome of their study include estimates for the probability of cholera elimination under several possible vaccination scenarios. Mimicking their approach, we define cholera elimination as an absence of cholera

infections for at least 52 consecutive weeks, and we provide forecasts under the following vaccination scenarios:

V0: No additional vaccines are administered.

V1: Vaccination limited to the departments of Centre and Artibonite, deployed over a two-year period.

V2: Vaccination limited to three departments: Artibonite, Centre, and Ouest deployed over a two-year period.

V3: Countrywide vaccination implemented over a five-year period.

V4: Countrywide vaccination implemented over a two-year period.

Simulations from probabilistic models (Models 1 and 3) represent possible trajectories of the dynamic system under the scientific assumptions of the models. Because Model 1 only accounts for national level disease dynamics, the pre-determined department-specific vaccination campaigns are carried out by assuming the vaccines are administered in one week to the same number of individuals that would have obtained vaccines if explicitly administered to the specific departments. We refer readers to Lee et al. (2020) and the accompanying supplement material for more details. Estimates of the probability of cholera elimination can therefore be obtained as the proportion of simulations from these models that result in cholera elimination. The results of these projections are summarized in Figs. 3.7.

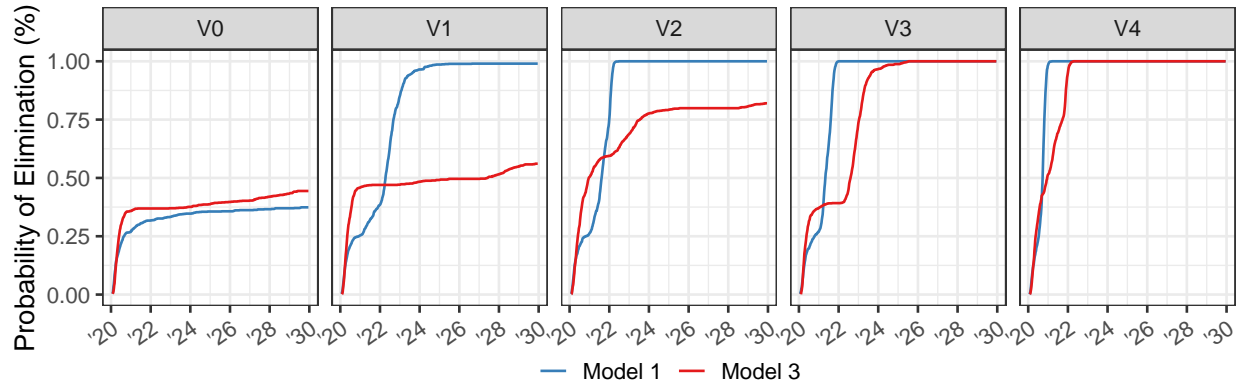


Figure 3.7: Probability of cholera elimination, defined as having zero cholera infectious for at least 52 consecutive weeks, based on 10 year simulations from calibrated versions of Models 1 and 3. Compare to Fig. 3A of Lee et al. (2020).

Probability of elimination estimates of this form are not meaningful for deterministic models, as the trajectory of these models only represent the mean behavior of the system

rather than individual potential outcomes. We therefore do not provide probability of elimination estimates under Model 2, but show trajectories under the various vaccination scenarios using this model (Fig. 3.4).

3.4 Discussion

The ongoing global COVID-19 pandemic has demonstrated how government policy may be affected by the inferences drawn from mathematical modeling (Saltelli et al., 2020). However, the development of credible models—which are supported by data and can provide quantitative insights into a dynamic system—remains a challenging task. In this article, we demonstrated opportunities available for raising the current standards of statistical inference for mathematical models of biological systems.

We presented methodology consistent with existing guidelines (Behrend et al., 2020) but going beyond standard practice. In particular, we showed the value of comparing the likelihood of fitted mechanistic models versus non-mechanistic benchmarks, a practice that has been previously advocated for (He et al., 2010) but was not done by any of the studies in our literature review. These comparisons, along with other likelihood based diagnostics, help identify specific limitations of proposed models. Diagnostic tools include likelihood profile methods, which help to assess parameter identifiability and enable the construction of confidence intervals for parameter estimates (Ionides et al., 2017; Simpson and Maclaren, 2023). When reaching conclusions, it is important to consider potential consequences of confounded variables and model misspecification.

Model diagnostics are a key tool for exposing unresolved model limitations and improving model fit. In our case study, we compared the three models from Lee et al. (2020) to statistical benchmarks, revealing areas for improvement. For example, comparisons of Model 3 to a benchmark revealed its inadequacy in accounting for the post-hurricane increase in transmission, leading to a beneficial model refinement. When a mechanistic model is competitive with statistical benchmarks, we have a license to begin critical evaluation of its causal implications. If a model falls far behind simple benchmarks, there is likely to be substantial limitation in the data analysis that should be identified and remedied. In our case study, the re-calibrated version of Model 1 outperformed its benchmark, so we proceeded to examine causal implications. When doing so, we found that the fitted model provides a causal description of the dynamic system that is consistent with known features of the system, such as the importance of rainfall as a driver of cholera infection. The congruency between causal implications of the model and our belief about the dynamic system, coupled with a strong quantitative description of observed data relative to a benchmark, provides

support for viewing the model as a plausible quantitative representation of the system under investigation.

When fitting a mechanistic model to a dynamic system, the complexity of the model warrants consideration. Mathematical models provide simplified representations of complex systems, with the simplicity serving both to facilitate scientific understanding and to enable statistical inference on unknown parameters. In our case study, employing deterministic dynamics in Model 2 was found to be an over-simplification by comparing model fit with benchmarks. Model 3 is distinct in that it is both stochastic and has a meta-population structure, making it challenging to draw likelihood-based inferences. In this paper, we demonstrated how this model class can be calibrated to incidence data using the innovative IBPF algorithm. One of only a few examples of fitting a nonlinear non-Gaussian meta-population model via maximum likelihood (Li et al., 2023; Ionides et al., 2022), this case study exemplifies the algorithm’s potential benefits and provides an example for future researchers on a possible approach to fitting a high-dimensional non-linear model.

Likelihood-based methods aid in determining an appropriate level of model complexity. Models fit to the same data can be compared using a criteria such as AIC. Nested model variations are particularly useful as they enable formal statistical testing of the nested features via likelihood ratio tests. Our case study demonstrated the examination of nested model features for all three models. Model 1 investigated a time-varying transmission rate; Model 2 assessed a phase-shift parameter in seasonal cholera peaks; Model 3 incorporated hurricane-related parameters.

Unmodeled features of a dynamic system can lead to spurious or misleading parameter estimates if the features substantially impact observed data. In deterministic models, features that cannot be explained by measurement error must be accounted for by the choice of parameters. For our case study, some of the parameter estimates for the deterministic Model 2 are implausible, such as the infinite immunity discussed above, and this may be explained by compensation for model misspecification. Incorporating demographic and environmental stochasticity into models can mitigate the impact of unmodeled features. Stochastic phenomena are not only arguably present in biological systems, but their inclusion in a model also allows observed data variations to be attributed to inherent uncertainty rather than to distorted parameter values. Models 1 and 3 suggest the presence of extra-demographic stochasticity (He et al., 2010; Stocks et al., 2020; Li et al., 2023), as evidenced by the confidence intervals for the corresponding parameter σ_{proc} (Appendix B.4).

If forecasts are an important component of a modeling task, the forecasts should be consistent with the available data, particularly at the most recently available time points. In our case study, we did this by simulating forward from the filtering distribution, as

this procedure conditions latent variables on the available data. This type of forecasting, however, is not directly available using a deterministic model, where future dynamics are fully determined by initial conditions and parameter values. This can result in over-confident model forecasts (King et al., 2015). Despite their limitations, deterministic models can offer valuable insights into dynamic systems (May, 2004). In Lee et al. (2020), the forecasts from the deterministic Model 2 were qualitatively more consistent with the observed disappearance of cholera than the stochastic models. In our case study, we found improvements to Models 1 and 3 that resulted in improved forecasts for these models.

In our case study, we found that additional attention to statistical details could have resulted in an enhanced statistical fit to the observed incidence data. This would have improved the accuracy of the policy guidance resulting from the study. We used the same data, models, and much of the same code used by Lee et al. (2020), but we arrived at drastically different conclusions. Specifically, each of the re-calibrated models predicted with moderate probability that cholera would disappear from Haiti. Although there have been new cases of cholera in Haiti, this conclusion aligns more with the prolonged absence of cholera cases from 2019-2022. We acknowledge the benefit of hindsight: our demonstration of a statistically principled route to obtain better-fitting models resulting in more robust insights does not rule out the possibility of discovering other models that fit well yet predict poorly.

Mechanistic models offer opportunities for understanding and controlling complex dynamic systems. This case study has investigated issues requiring attention when applying powerful new statistical techniques that can enable statistically efficient inference for a general class of partially observed Markov process models. Researchers should ensure that intensive numerical calculations are adequately executed. Using benchmarks and alternative model specifications to assess statistical goodness-of-fit should also be common practice. Once a model has been adequately calibrated to data, care is required to assess what causal conclusions can properly be inferred given the possibility of alternative explanations consistent with the data. Studies that combine model development with thoughtful data analysis, supported by a high standard of reproducibility, build knowledge about the system under investigation. Cautionary warnings about the difficulties inherent in understanding complex systems (Saltelli et al., 2020; Ioannidis et al., 2020; Ganusov, 2016) should motivate us to follow best practices in data analysis, rather than avoiding the challenge.

3.4.1 Reproducibility and Extendability

Lee et al. (2020) published their code and data online, and this reproducibility facilitated our work. Robust data analysis requires not only reproducibility but also extendability: if

one wishes to try new model variations, or new approaches to fitting the existing models, or plotting the results in a different way, this should not be excessively burdensome. Scientific results are only trustworthy so far as they can be critically questioned, and an extendable analysis should facilitate such examination (Gentleman and Temple Lang, 2007).

We provide a strong form of reproducibility, as well as extendability, by developing our analysis in the context of a software package, `haitipkg`, written in the R language (R Core Team, 2022). Using a software package mechanism supports documentation, standardization and portability that promote extendability. In the terminology of Gentleman and Temple Lang (2007), the source code for this article is a *dynamic document* combining code chunks with text. In addition to reproducing the article, the code can be extended to examine alternative analysis to that presented. The dynamic document, together with the R packages, form a *compendium*, defined by Gentleman and Temple Lang (2007) as a distributable and executable unit which combines data, text and auxiliary software (the latter meaning code written to run in a general-purpose, portable programming environment, which in this case is R).

CHAPTER 4

Conclusion

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

APPENDIX A

Appendix for Chapter 2

A.1 Comparing ARIMA algorithm to Python software

Although Algorithm 1 is language agnostic, our implementation is in R. Consequently, all simulation studies for this article use R’s `stats::arima` function for baseline comparisons. To demonstrate applicability to other software environments, we briefly compare model likelihoods fit using Python’s `statsmodels.tsa` module against our implementation of Algorithm 1. We generate 100 unique Gaussian ARMA(2,1) models and datasets (each with $n = 100$) where R’s `stats::arima` provides sub-optimal estimates. These observations are used to fit ARMA(2,1) models in Python. Despite these datasets being chosen for sub-optimal results in R, the log-likelihoods in both R and Python are roughly equivalent—a result of both software packages using the same general approach to fitting model parameters.

Our implementation of Algorithm 1 resulted in higher log-likelihoods for 97 out of the 100 datasets compared to Python. The log-likelihood deficiencies for the remaining three datasets were all smaller than $\epsilon = 10^{-5}$, which is smaller than the tolerance level to be considered as an improvement in our other simulation studies. While potentially insignificant, these differences can be eliminated by increasing the number of parameter initializations and the convergence criteria of our algorithm. Directly implementing our algorithm in Python would further eliminate the possibility of these discrepancies entirely.

A.2 Uniform Sampling

A common approach to optimizing a non-convex loss function is to perform the optimization routine with distinct parameter initializations. For ARMA models, picking suitable initialization is a challenging problem that we address with Algorithm 1. An alternative approach would involve sampling each coefficient independently. To see why an independent sampling scheme is not used, consider an AR(2) model. An initialization with parameters

$(\phi_1, \phi_2) = (1.1, 0.1)$ is not a valid initialization—as the polynomial roots lie outside the complex unit circle—whereas $(\phi_1, \phi_2) = (1.1, -0.2)$ is perfectly acceptable. Our algorithm accounts for the complex relationship between model parameters when obtaining random initializations.

To visualize why an independent sampling scheme is not used, consider sampling parameters from a $\text{Uniform}(-1, 1)$ distribution. In Fig A.1, we plot inverted roots that are a result of sampling from this distribution for AR(2) and AR(3) models. The figure illustrates that a significant percentage of uniformly sampled parameter initializations lie outside the accepted region, and, critically, the entire region of possible initializations is not well covered by uniform sampling. Picking a uniform distribution with different bounds—or any other independent sampling distribution—results in similar problems. In order to uniformly sample from the possibly regions, it is necessary to account for the geometry of parameter space, a problem solved by Algorithm 1.

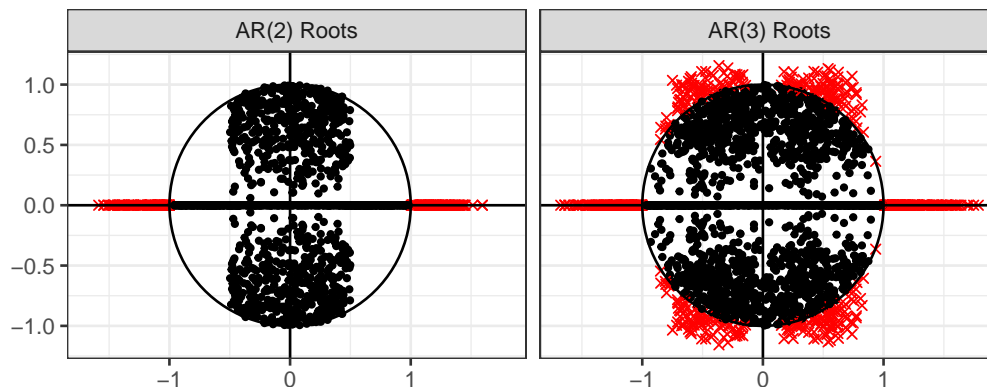


Figure A.1: Inverted roots of 1000 samples of AR(2) and AR(3) coefficients sampled independently from a $U(-1, 1)$ distribution. The red “x”s are points that lie outside the accepted region, which represents 12.2% of the AR(2) coefficients and 22.8% of the AR(3) coefficients. Increasing the width of the uniform sampling distribution results in a larger fraction outside the accepted region, and decreasing the width results in worse coverage of the range of acceptable parameter values.

APPENDIX B

Appendix for Chapter 3

B.1 Model Diagrams

Each of the dynamic models considered in this manuscript can be fully described using the model descriptions in the manuscript, coupled with the additional information described in Sections 2 and 3 of this supplement. Despite this, diagrams of dynamic systems are often helpful to understand the equations. In this section, we give three diagrams representing Models 1–3, respectively. Because the models are defined by their mathematical equations and numeric implementation, these diagrams are not unique visual representations of the model. Alternative representations that may be helpful in understanding the models explored in this paper are provided in the supplement material of Lee et al. (2020).

B.2 Markov chain and differential equation interpretations of compartment flow rates

In the Materials and methods Section of the main article, we define compartment models in terms of their flow rates. For a discrete population model, these rates define a Markov chain. For a continuous and deterministic model, the rates define a system of ordinary differential equations. Here, we add additional details to clarify the mapping from a collection of rate functions to a fully specified process. Our treatment follows Bretó et al. (2009).

A general compartment model is a vector-valued process $X(t) = (X_1(t), \dots, X_c(t))$ denoting the (integer or real-valued) counts in each of c compartments, where t is any continuous value in the interval $[t_0, \infty)$ for some real valued starting time t_0 . The compartments may also have names, but to set up general notation we simply refer to them by their numerical index. The basic characteristic of a compartment model is that $X(t)$ can be written in terms of the flows $N_{ij}(t)$ from i to j . A flow into compartment i from outside the system is denoted by $N_{\bullet i}$, and a flow out of the system from compartment i is denoted by $N_{i\bullet}$. We call \bullet a

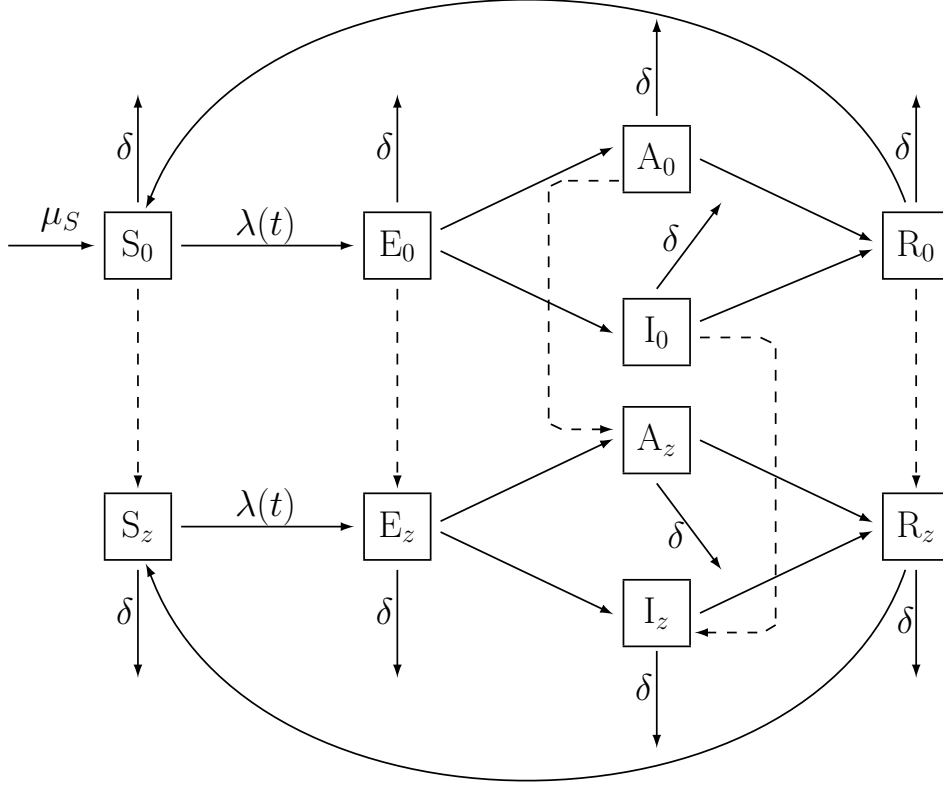


Figure B.1: A flow diagram for the SEAIR model.

source/sink compartment, though it is an irregular compartment since $X_{\bullet}(t)$ is not defined. These flows are required to satisfy a “conservation of mass” identity:

$$X_i(t) = X_i(t_0) + N_{\bullet i}(t) - N_{i\bullet}(t) + \sum_{j \neq i} N_{ji}(t) - \sum_{j \neq i} N_{ij}(t). \quad (\text{B.1})$$

Each flow $N_{ij}(t)$ is associated with a rate function $\mu_{ij} = \mu_{ij}(t, X(t))$, where we include the possibility that i or j takes value \bullet .

There are different ways to use a collection of rate functions to build a fully specified model. We proceed to describe the ones we use in this paper: via a system of ordinary differential equations (Sec. B.2.1), a simple Markov counting system (Sec. B.2.2), and an over-dispersed Markov counting system (Sec. B.2.3). Other representations include stochastic

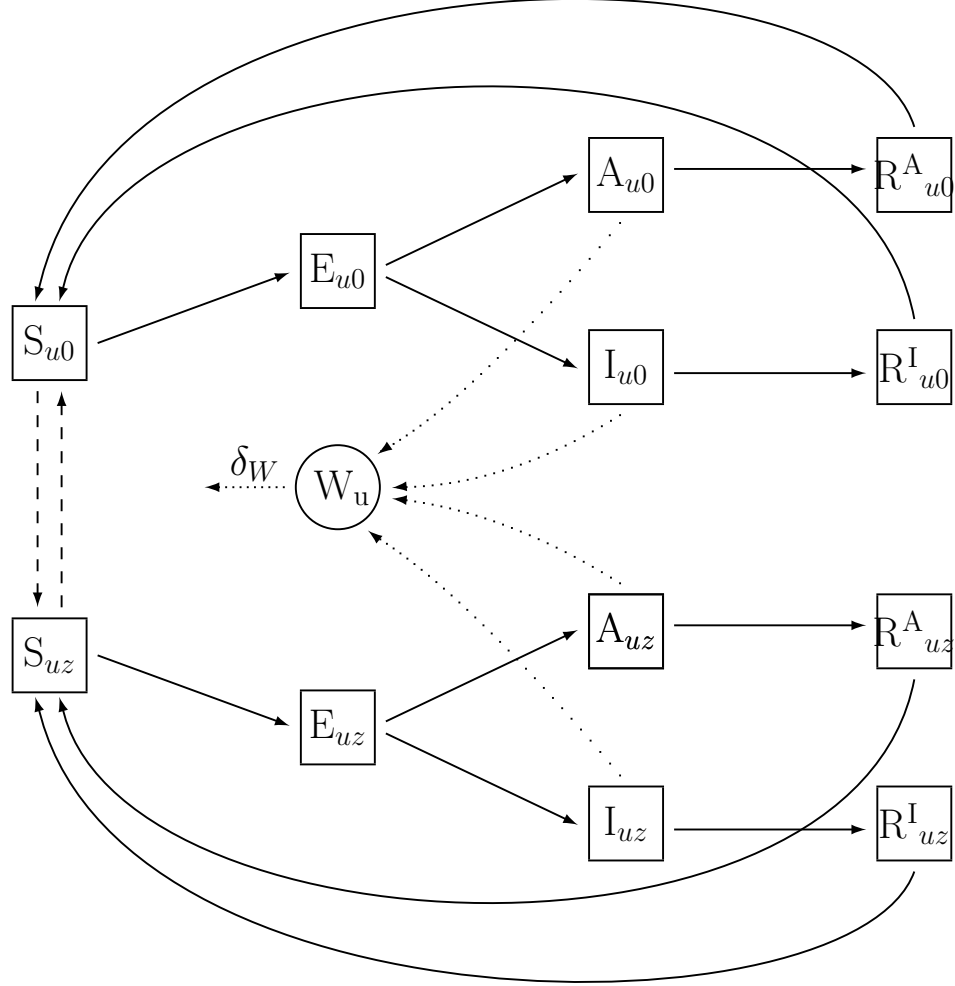


Figure B.2: A flow diagram for the SEAIR model 2. This is a constant population model, there are no births/deaths. Vaccinations are assumed to only be given to susceptible individuals, and vaccine immunity wanes only with susceptible vaccinated individuals.

differential equations driven by Gaussian noise or Gamma noise Bhadra et al. (2011).

B.2.1 Ordinary differential equation (ODE) interpretation

A basic deterministic specification is

$$dN_{ij}/dt = \mu_{ij}(t, X(t))X_i(t), \quad i \in 1:c, \quad j \in 1:c \cup \{\bullet\}, \quad i \neq j, \quad (\text{B.2})$$

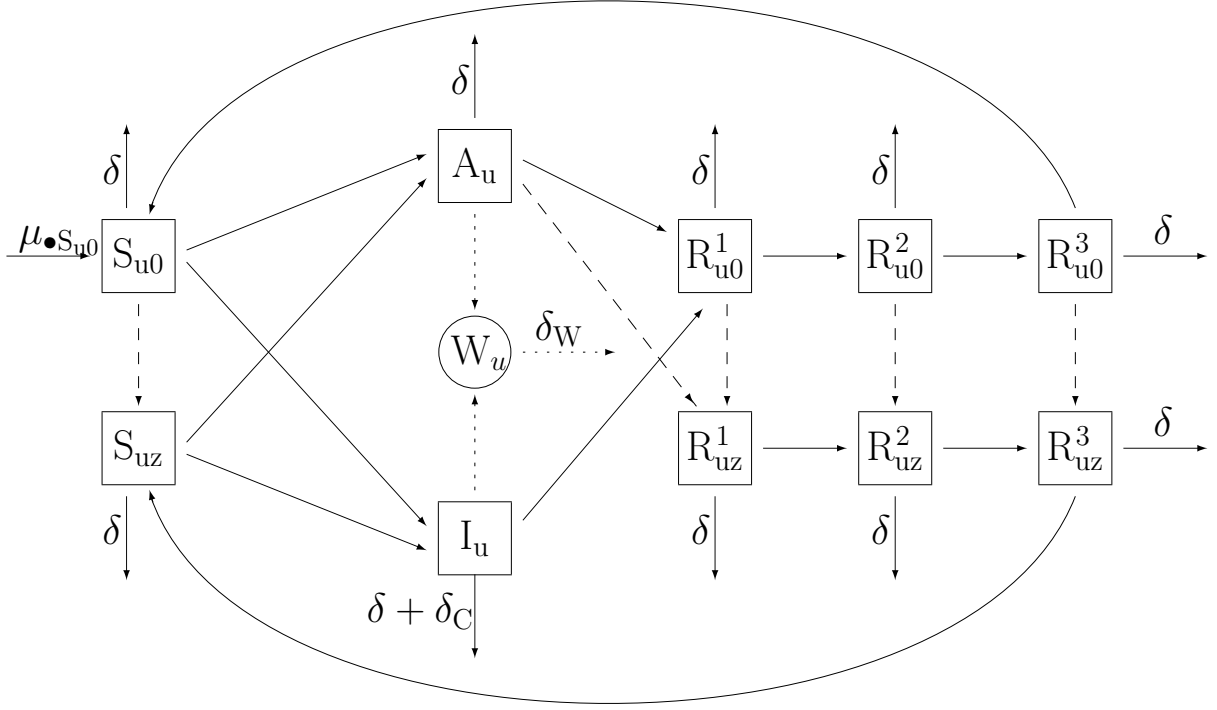


Figure B.3: A flow diagram for the SAIR model 3. This model assumes a constant population while also including a mechanism for births/deaths; all deaths are balanced by births into the unvaccinated susceptible compartment, so the birth rate $\mu_{\bullet} S_{u0}$ corresponds to the sum total deaths from the remaining compartments. The model assumes that symptomatic individuals will not be vaccinated, hence no vaccination arrow exiting the I_{u0} compartment.

where $\mu_{ij}(t, X(t))$ is called a per-capita rate or a unit rate. Flows into the system require special treatment since $X_i(t)$ in (B.2) is not defined for $i = \bullet$. Instead, we specify

$$dN_{\bullet i}/dt = \mu_{\bullet i}(t, X(t)). \quad (\text{B.3})$$

This is the the interpretation and implementation used for Model 2 in our study.

B.2.2 Simple Markov counting system interpretation

A continuous time Markov chain can be specified via its infinitesimal transition probabilities. A basic approach to this is to define

$$\mathbb{P}[N_{ij}(t + \delta) - N_{ij}(t) = 0 \mid X(t)] = 1 - \delta\mu_{ij}(t, X(t))X_i(t) + o(\delta), \quad (\text{B.4})$$

$$\mathbb{P}[N_{ij}(t + \delta) - N_{ij}(t) = 1 \mid X(t)] = \delta\mu_{ij}(t, X(t))X_i(t) + o(\delta), \quad (\text{B.5})$$

for $i \in 1:c$ and $j \in 1:c \cup \{\bullet\}$ with $i \neq j$. As with the ODE case, we need special attention for flows into the system, and we define

$$\mathbb{P}[N_{\bullet i}(t + \delta) - N_{\bullet i}(t) = 0 \mid X(t)] = 1 - \delta\mu_{\bullet i}(t, X(t)) + o(\delta), \quad (\text{B.6})$$

$$\mathbb{P}[N_{\bullet i}(t + \delta) - N_{\bullet i}(t) = 1 \mid X(t)] = \delta\mu_{\bullet i}(t, X(t)) + o(\delta). \quad (\text{B.7})$$

Together with the initial conditions $X(0)$, equations (B.4)–(B.7) define a Markov chain. Each flow is a simple counting process, meaning a non-decreasing integer-valued process that only has jumps of size one. We therefore call the Markov chain a simple Markov counting system (SMCS). The infinitesimal mean of every flow is equal to its infinitesimal variance Bretó and Ionides (2011) and so an SMCS is called equidispersed. We note that the special case of Model 1 used by Lee et al. (2020) (with $\sigma_{\text{proc}} = 0$) is an SMCS. To permit more general mean-variance relationships for a Markov counting system, we must permit jumps of size greater than one. The utility of over-dispersed models, where the infinitesimal variance of the flow exceeds the infinitesimal mean, has become widely recognized Stocks et al. (2020); He et al. (2010).

B.2.3 Overdispersed Markov counting system interpretation

Including white noise in the rate function enables the possibility of an over-dispersed Markov counting system Bretó and Ionides (2011); Bretó et al. (2009); He et al. (2010). Since rates should be non-negative, Gaussian noise is not appropriate and gamma noise is a convenient option that has found various applications Romero-Severson et al. (2015); Subramanian et al. (2020). Specifically, we consider a model given by

$$\mu_{ij}(t, X(t)) = \bar{\mu}_{ij}(t, X(t)) d\Gamma_{ij}(t)/dt, \quad (\text{B.8})$$

where $\Gamma_{ij}(t)$ is a stochastic process having independent gamma distributed increments, with

$$\mathbb{E}[\Gamma_{ij}(t)] = t, \quad \text{Var}[\Gamma_{ij}(t)] = \sigma_{ij}^2 t. \quad (\text{B.9})$$

Formally interpreting the meaning of (B.8) is not trivial, and we do so by constructing a Markov process $X(t)$ as the limit of the Euler scheme described in Section B.3, below. Therefore, the numerical scheme in Sec. B.3 can be taken as a definition of the meaning of (B.8). The Markov chain defined by the limit of this Euler scheme as the step size decreases is an over-dispersed Markov counting system, with the possibility of instantaneous jumps of size greater than one Bretó and Ionides (2011).

B.3 Numerical solutions to compartment models

Models may be fitted and their implications assessed via numerical solutions (i.e., simulations) from the model equations. All the analyses we consider have this simulation-based property, known as plug-and-play or equation-free or likelihood-free. The numerical solutions to the model are arguably of more direct scientific interest than the exact solutions to the postulated equations. For ODE models, numerical methods are well studied and a standard numerical solution package such as `deSolve` in R is adequate for many purposes. For SMCS and ODMCS models, exact schemes are feasible when the number of events is small, which may be the case for small populations. However, for applicability to larger populations, we use instead the following Euler scheme. Write δ for an Euler time step, and ΔN_{ij} for the numerical approximation to $N_{ij}(t + \delta) - N_{ij}(t)$ given $X(t)$. For each i and j in $1:c \cup \{\bullet\}$ with $i \neq j$, we draw independent Gamma distributed noise increments with mean δ and variance $\sigma_{ij}^2 \delta$, denoted using a mean-variance parameterization of the gamma distribution as

$$\Delta \Gamma_{ij} \sim \text{gamma}(\delta, \sigma_{ij}^2 \delta). \quad (\text{B.10})$$

In the case of an SMCS model, $\sigma_{ij} = 0$ for all i and j , so we have $\Delta \Gamma_{ij} = \delta$. Then, for $i \neq \bullet$ and $j \neq i$, and writing

$$\mu_{ij} = \bar{\mu}_{ij}(t, X(t)) \Delta \Gamma_{ij} / \delta, \quad (\text{B.11})$$

we calculate transition probabilities

$$p_{ij} = \exp \left\{ - \sum_{k \in 1:c \cup \{\bullet\}} \mu_{ik} \delta \right\} \frac{\mu_{ij}}{\sum_{k \in 1:c \cup \{\bullet\}} \mu_{ik}}, \quad (\text{B.12})$$

$$p_{ii} = 1 - \sum_{j \neq i} p_{ij}. \quad (\text{B.13})$$

These probabilities correspond to competing hazards for every individual in compartment i to transition to some compartment j , interpreting $j = i$ to mean that the individual remains in i . Then, $(\Delta N_{i1}, \dots, \Delta N_{ic}, \Delta N_{i\bullet})$ has the multinomial distribution where $X_i(t)$ individuals

are allocated independently to $1 : c \cup \{\bullet\}$ with probabilities given by (B.12) and (B.13). We use the `reulermultinom` function in the `pomp` package to draw from this multinomial distribution.

Different treatments of demographic flows—such as birth, death, immigration and emigration—are possible. For the case $i = \bullet$, the treatment used by Model 1 is to set

$$\Delta N_{\bullet j} \sim \text{poisson}(\mu_{\bullet j} \delta), \quad (\text{B.14})$$

an independent Poisson random variable with mean $\mu_{\bullet j} \delta$.

Models 2 and 3 used an alternative approach, balancing the total number of flows in and out of the compartment, i.e., $\sum_i N_{\bullet i}(t) = \sum_i N_{i\bullet}(t)$, in order to make the model consistent with the known total population. In this case, we formally model the death rate as a rate of returning to the susceptible class S , and use external transitions from \bullet into S to describe only net population increase.

B.4 Confidence Intervals for Model Parameters

In this section we provide confidence intervals for all model parameters, excluding those that take unique values for each spatial unit. For each model and parameter, we use principles of profile likelihood to obtain confidence intervals (Pawitan, 2001). Due to the non-linear and stochastic nature of Models 1 and 3, exact evaluations of the profile log-likelihood are difficult to obtain. Instead, the log-likelihood at each point of the profile is estimated using via Monte-Carlo based particle filter methods. We therefore obtain confidence intervals for the parameters of Model 1 and Model 3 using the Monte Carlo adjust profile (MCAP) algorithm (Ionides et al., 2017).

Profile confidence intervals for nonlinear POMP models are require a large number of computations. In the Model 1 and Model 3 subsections, we mention the total computational expense of each profile log-likelihood evaluation. Each subsection also provide figures that show the curvature of the profile log-likelihood near the MLE (Figures B.4–B.6). In these figures, the parameter values are shown on the transformed scale in which the profile was calculated.

B.4.1 Model 1 parameters

Parameter estimates for Model 1, along with the MCAP confidence intervals for the estimate, are given in Table B.1. Figure B.4 displays the Monte Carlo evaluations of the profile likelihood values, obtained using a particle filter. The total computational burden of this

profile likelihood search was 3631 hours, which was computed in parallel using 9675 separate jobs via the `batchtools` R package Lang et al. (2017).

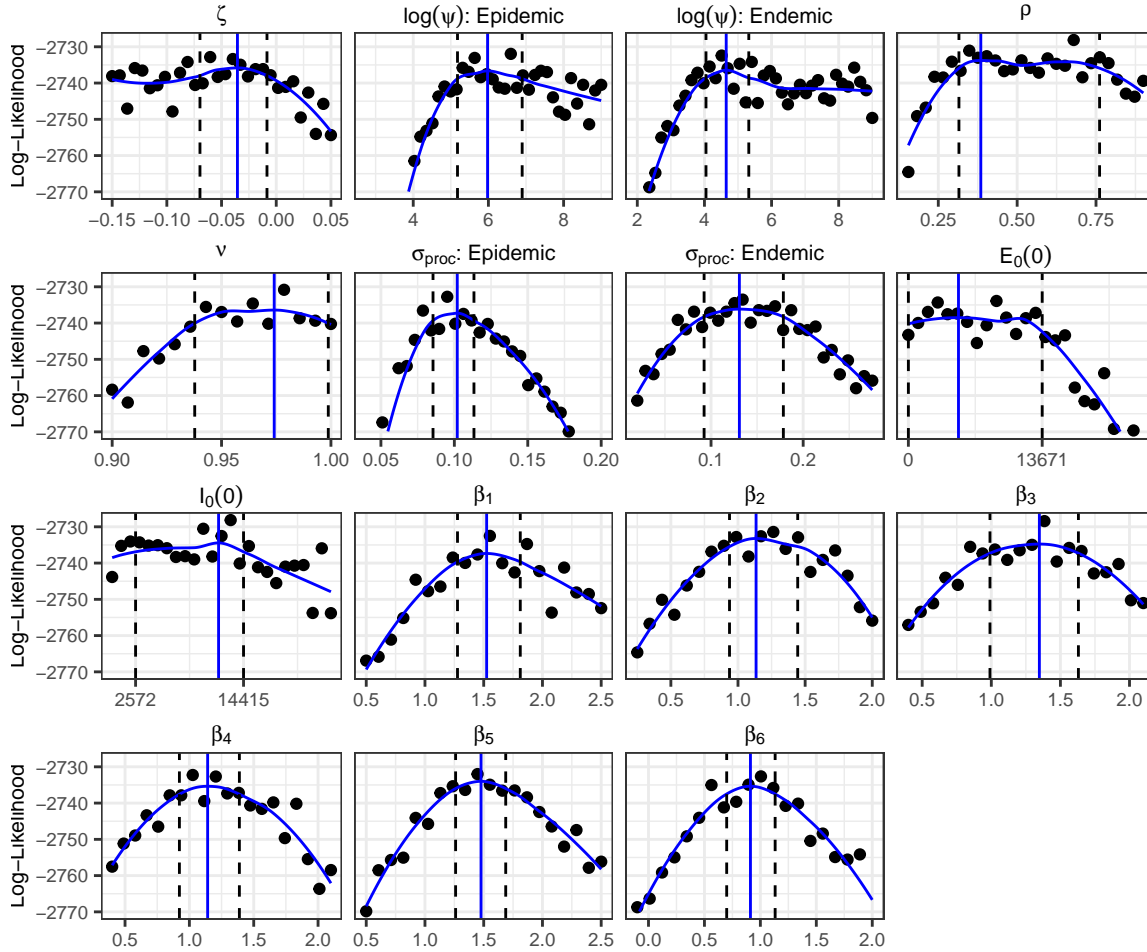


Figure B.4: MCAP confidence intervals for Model 1 parameters. The vertical blue line indicates the smoothed MLE.

B.4.2 Model 2 parameters

Parameter estimates for Model 2, along with the profile likelihood confidence intervals for each estimate, are given in Table B.2. Figure B.5 displays the profile log-likelihood curve near the MLE. In Table B.2, the confidence interval for μ_{RS}^{-1} , the duration of natural immunity due to cholera infection, is arbitrarily large (going to infinity). This is possible because the parameter that was estimated was μ_{RS} , and the true MLE for this parameter is zero (see Figure B.5). This suggests that the fitted model favors a regime where reinfection events are not possible. Similarly, the MLE for the parameter β , which controls the amount of cholera transmission from human to human, is zero. Because Model 2 fails to describe the

Table B.1: Model 1 parameter estimates and their corresponding confidence intervals, obtained via the MCAP algorithm.

Mechanism	Parameter	MLE	95% Confidence Interval
Seasonality	ζ	-0.036	(-0.070, -0.008)
Seasonality	β_1	1.417	(1.277, 1.811)
Seasonality	β_2	1.169	(0.937, 1.445)
Seasonality	β_3	1.136	(0.990, 1.630)
Seasonality	β_4	1.140	(0.922, 1.389)
Seasonality	β_5	1.401	(1.261, 1.687)
Seasonality	β_6	0.988	(0.699, 1.132)
Observation Variance	ψ : Epi	279.147	(177.226, 990.191)
Observation Variance	ψ : End	78.326	(57.171, 204.654)
Reporting Rate	ρ	0.679	(0.315, 0.761)
Mixing Exponent	ν	0.978	(0.938, 0.999)
Process noise (wk ^{1/2})	σ_{proc} : Epi	0.092	(0.085, 0.113)
Process noise (wk ^{1/2})	σ_{proc} : End	0.118	(0.092, 0.179)
Initial Values	$I_0(0)$	7298	(2572, 1.4415×10^4)
Initial Values	$E_0(0)$	350	(1, 1.3671×10^4)

incidence data as well as a simple statistical benchmark, we must be careful to not interpret these results as evidence that reinfections and human-to-human infection events do not occur. Instead, we may consider this as additional evidence of model misspecification.

Table B.2: Model 2 parameter estimates and their corresponding confidence intervals, obtained via profile likelihood.

Mechanism	Parameter	MLE	95% Confidence Interval
Human to water shedding (wk ⁻¹)	μ_W	179.2	(144.6, 229.4)
Water to Human Infection (yr ⁻¹)	β_W	1.098	(1.067, 1.128)
Observation Variance	ψ	1.319	(1.291, 1.347)
Seasonality	ϕ	0.974	(7.127, 7.381)
Human to Human Infection (yr ⁻¹)	β	$5.97 \times 10^{-15*}$	$[0, 2.3 \times 10^{-6}]$
Immunity (yr)	μ_{RS}^{-1}	$1.4 \times 10^{11*}$	(1410, inf)

* As evident in Figure B.5, the true MLE for these parameters is 0 and ∞ , respectively; this value could not be obtained numerically due to the parameter transformation applied to the parameter for the model fitting processes.

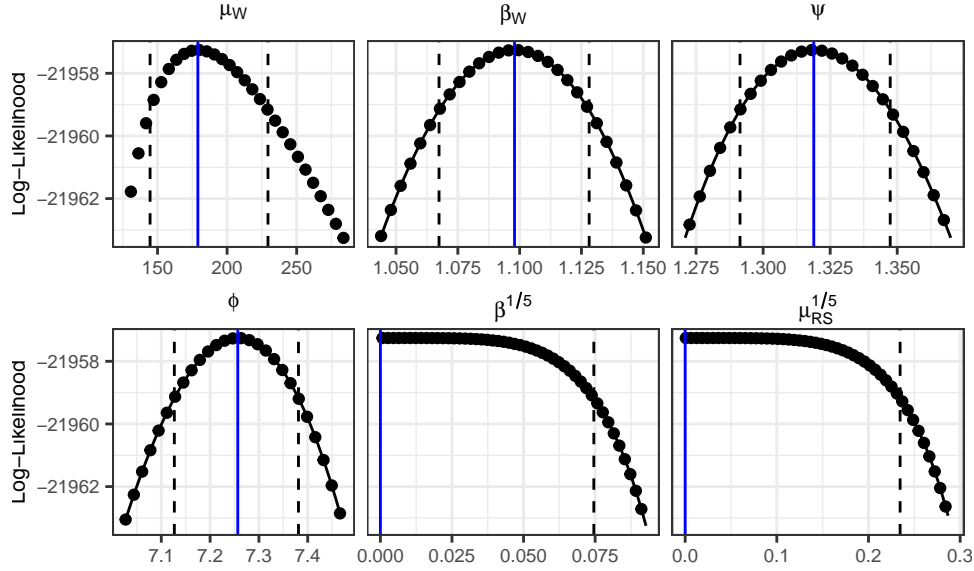


Figure B.5: MCAP confidence intervals for Model 2 parameters. The vertical blue line indicates the MLE.

B.4.3 Model 3 parameters

Parameter estimates for Model 3, along with the MCAP confidence intervals for the estimate, are given in Table B.3. Figure B.6 displays the Monte Carlo evaluations of the profile likelihood values, obtained using a particle filter. The total computational burden of this profile likelihood search was 28938 hours, which was computed in parallel using 7568 separate jobs via the `batchtools` R package Lang et al. (2017).

Table B.3: Model 3 parameter estimates and their corresponding confidence intervals, obtained via the MCAP algorithm.

Mechanism	Parameter	MLE	95% Confidence Interval
Process Noise ($\text{wk}^{1/2}$)	σ_{proc}	0.218	(0.203, 0.230)
Water Survival (wk)	δ_W^{-1}	0.108	(0.087, 0.110)
Human to Water Shedding $\frac{\text{km}^2}{\text{wk}}$	μ_W	9.77×10^{-7}	$(8.64 \times 10^{-7}, 1.25 \times 10^{-6})$
Asymptomatic Shedding	ϵ_W	0.008	(0.0, 0.095)
Seasonality	a	1.000	(0.637, 1.432)
Seasonality	r	0.780	(0.498, 1.041)
Reporting Rate	ρ	0.983	(0.789, 1.000)
Observation Variance	ψ	88.578	(66.034, 132.563)

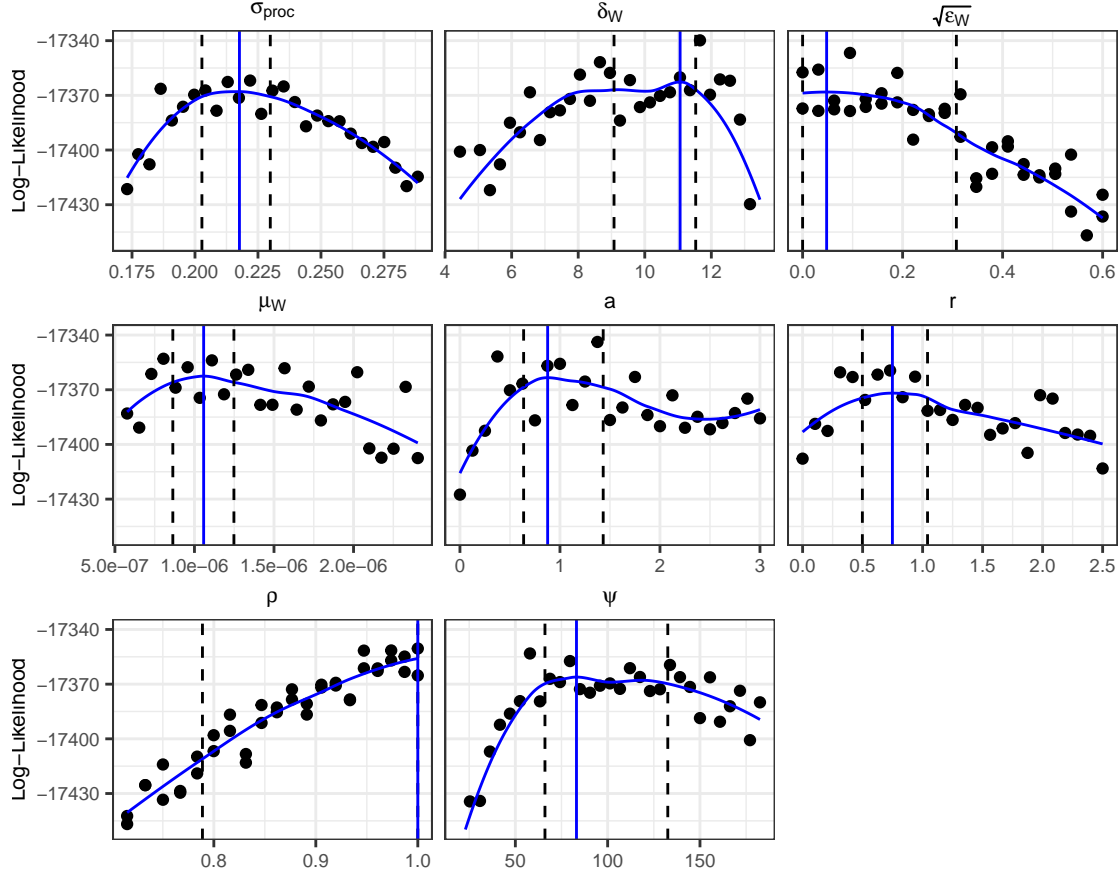


Figure B.6: MCAP confidence intervals for Model 3 parameters. The vertical blue line indicates the smoothed MLE.

B.5 Replication of Lee et al. (2020) results

In this article we claimed that we were able to obtain better fits to the observed data using the same models that were proposed by Lee et al. (2020). Along with visual comparisons to the data, this claim was supported by comparing likelihoods and AIC values in Table 2 in the manuscript. Because model likelihoods were not provided by Lee et al. (2020), it is necessary to replicate these models in order to obtain likelihood estimates. Here we would like to thank the authors of Lee et al. (2020), who provided detailed descriptions of their models, which enabled us to build on their work. In the following subsections, we use our R package `haitipkg` to reproduce some of the results of Lee et al. (2020). This reproduction allows us to estimate the likelihoods of the Lee et al. (2020) version of Models 1–3, and also provides a demonstration of the importance and usefulness of reproducible research.

B.5.1 Model 1 Replication

The model was implemented by a team at Johns Hopkins Bloomberg School of Public Health (hereafter referred to as the Model 1 authors) in the R programming language using the `pomp` package (King et al., 2016). Original source code is publicly available with DOI: 10.5281/zenodo.3360991. The final results reported by the Model 1 authors were obtained by using several different parameter sets rather than a single point estimate. According to the supplement materials, this was because model realizations from a single parameter set retained substantial variability, but multiple realizations from a collection of parameter sets resulted in a reasonable visual fit to the data. We are also inclined to believe that the use of multiple parameter values was in part intended to account for parameter uncertainty—the importance of which was discussed in the main text—an effort by the Model 1 authors that we applaud. Simulations from each of the parameter sets, however, were treated with equal importance when being used to diagnose the model fit and make inference on the system. This is problematic given Figures S8 and S9 of the supplement material, which suggest that some parameter sets that were used for inference may have been several hundred units of log-likelihood lower than other parameter sets that were simultaneously used to make forecasts. Such a large difference in log-likelihoods is well beyond the threshold of statistical uncertainty determined by Wilks’ theorem, resulting in the equal use of statistically inferior parameter sets in order to make forecasts and conduct inference on the system.

To fully reproduce the results of the Model 1 authors, it is necessary to use the exact same set of model parameters that were originally used to obtain the results presented by Lee et al. (2020). Because these parameter sets were not made publicly available, we relied on the source code provided by the Model 1 authors to approximately recreate the parameter set. Due to software updates since the publication of the source code, we were unable to produce the exact same set of parameters. Running the publicly available source code, however, resulted in a set of parameters that are visually similar to those used by the Model 1 authors (See Figures B.7 and B.8). Furthermore, simulations using the set of parameters produced by the source code appear practically equivalent to those displayed by Lee et al. (2020) (See Figure B.9).

Because the model forecasts provided by Lee et al. (2020) come from various sets of parameters—which each correspond to a unique log-likelihood value—it is not obvious how one would obtain an estimate for the log-likelihood of the model that was used for simulations by the Model 1 authors. One approach could be to calculate the logarithm of the weighted mean of the likelihoods for each parameter sets used to obtain the forecasts, where the weights are proportional to the number of times the parameter set was used. However, in an effort to not underestimate the likelihood of the model of the Model 1 authors, we report the

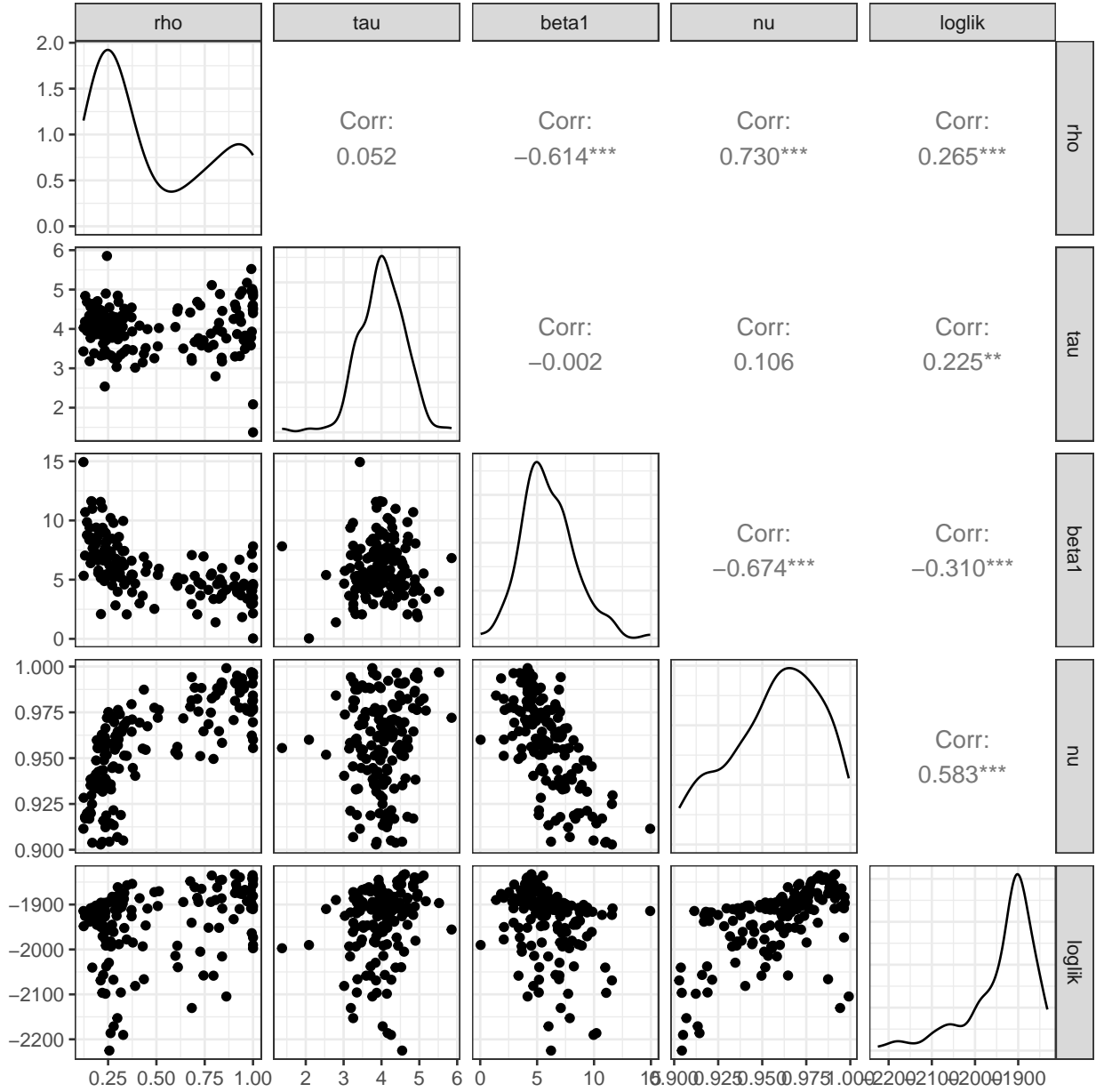


Figure B.7: Bivariate distributions of parameter estimates after fitting epidemic phase of the Model 1 following the procedure described by Lee et al. (2020). Compare to Figure S8 in the supplement of Lee et al. (2020).

estimated log-likelihood as the log-likelihood value corresponding to the parameter set with the largest likelihood value, even though the majority of simulations were obtained using parameter sets with lower likelihood values. In this sense, we consider the log-likelihood reported in Table 1 of the main text to be an upper-bound of the log-likelihood of the model used by Lee et al. (2020). For each parameter set, the log-likelihood was estimated using a

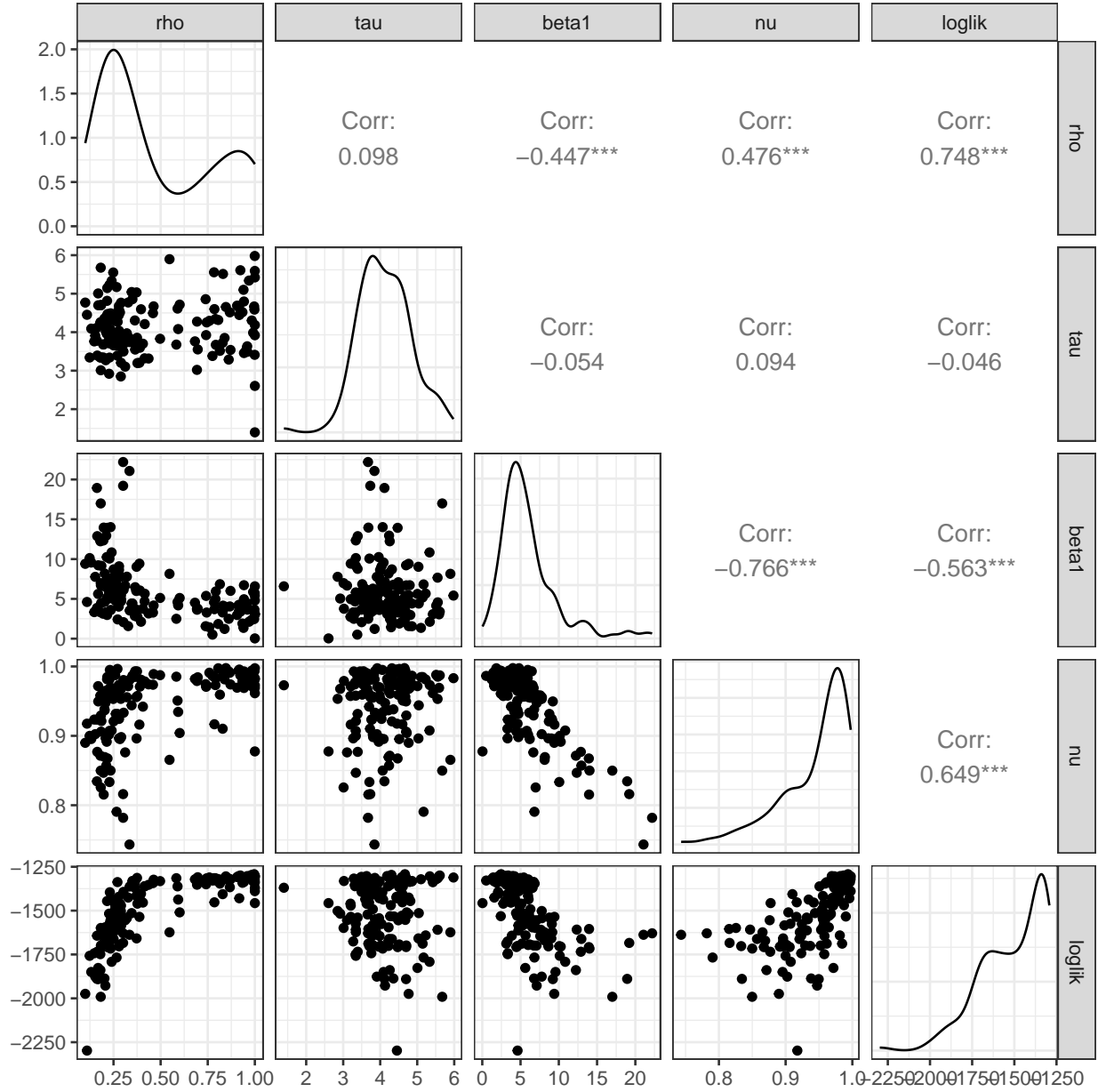


Figure B.8: Bivariate distributions of parameter estimates after fitting endemic phase of the Model 1 following the procedure described by Lee et al. (2020). Compare to Figure S9 in the supplement of Lee et al. (2020).

particle filter, implemented as the `pfilter` function in the `pomp` package.

B.5.2 Model 2 Replication

Model 2 was developed by a team that consisted of members from the Fred Hutchinson Cancer Research Center and the University of Florida (hereafter referred to as the Model 2 authors).

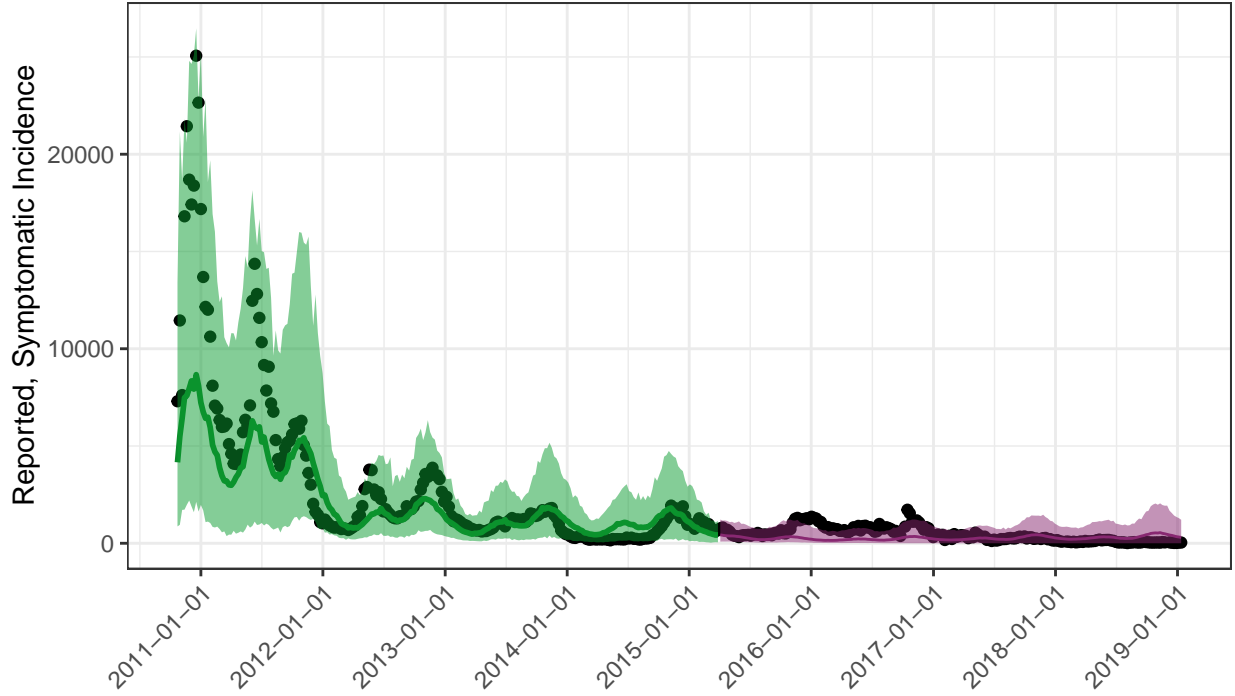


Figure B.9: Simulations from Model 1 using parameter sets that were generated by running source code provided by Lee et al. (2020). Compare to Figure S7 in the supplement of Lee et al. (2020). The upper bound for the likelihood of this model is -3031.

While Model 2 is the only deterministic model we considered in our analysis, it contains perhaps the most complex descriptions of cholera in Haiti: Model 2 accounts for movement between spatial units; human-to-human and environment-to-human cholera infections; and transfer of water between spatial units based on elevation charts and river flows.

The source code that the Model 2 authors used to generate their results was written in the `Python` programming language, and is publicly available at [10.5281/zenodo.3360857](https://zenodo.org/record/3360857) and its accompanying GitHub repository <https://github.com/lulelita/HaitiCholeraMultiModelingProject>. In order to perform our analysis in a unified framework, we re-implemented this model in the `R` programming language using the `spatPomp` package (Asfaw et al., 2024), which facilitates the creation of meta-population models. We note that the travel and water matrices used to implement the complex dynamics in Model 2 are not available in either the Zenodo archive or the GitHub repository; instead, we obtained these matrices via personal correspondence with the Model 2 authors. Using these matrices, and the point estimates for model parameters provided by Lee et al. (2020), we created trajectories of the cholera dynamics and compared this to available data. These trajectories, shown in Figure B.10, are very similar to the trajectories shown in

Figure S15 of the supplement of Lee et al. (2020).

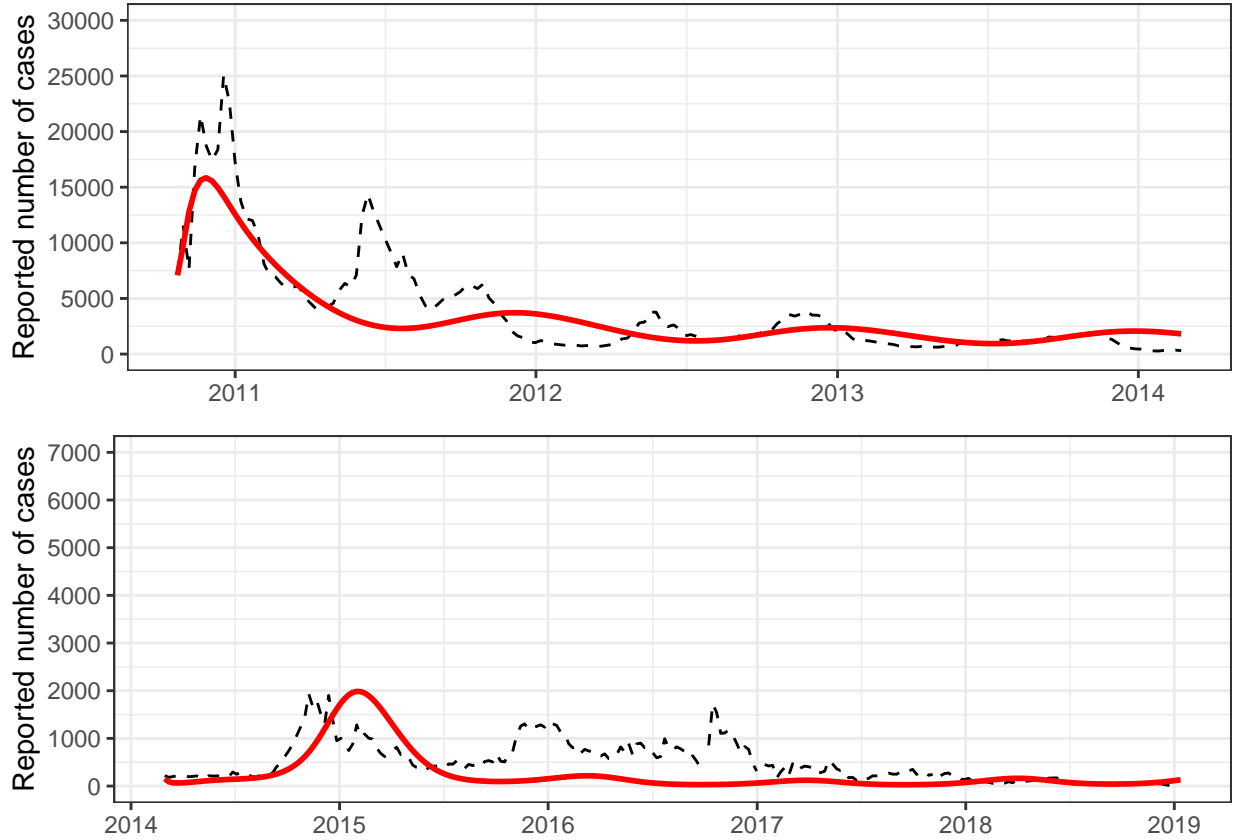


Figure B.10: Model 2 trajectories using the `haitipkg`. Compare to Figure S15 in the supplement of Lee et al. (2020).

There are minor differences between Figure B.10 and Figure S15 of Lee et al. (2020). While the discrepancy appears minor, the deterministic nature of Model 2 implies that an exact replication of model trajectories should be possible. In this case, these discrepancies may possibly be attributed to implementing the model and plotting the model trajectory in two different programming languages. Another potential explanation for the discrepancy is that the parameters that we used are only approximately the same as those used by Lee et al. (2020). For example, the parameters β , β_W had reported values of 9.9×10^{-7} and 4.03×10^{-2} , respectively (Table S13 of the supplement material of Lee et al. (2020)), but were actually fit to data and therefore likely these values have been rounded. Additionally, our implementation of Model 2 used a time scale of years and many of the parameters were reported on a weekly scale, so small differences may result due to unit conversions. The collective effect of these small differences in model parameters likely will result in small

differences in model trajectories.

Some additional concerns about being able to accurately replicate the results of Lee et al. (2020) are valid. Details about the measurement models and how latent states were initialized for the epidemic model were not provided by Lee et al. (2020) and therefore these details must be inferred by looking at the provided source code. According to repository comments, the files `fitInPieces3paramsCleanMay2019Public.py` and `fitInPiecesMuWithFracSusFixedAllInfectionsPublic.py` were used to fit the epidemic and endemic phases of the model respectively, although it is apparent that these exact files were not used to obtain the reported results since the files contain some variable-naming errors that make it impossible to run the files without making modifications ¹. The inability to replicate the results by Lee et al. (2020) by running the provided source code makes checking whether or not a our numeric implementation faithfully represents their results very difficult. Additionally, the script that was said to been used to obtain the results reported by Lee et al. (2020) appears to use a different measurement model than what was described in the supplemental material, again making it difficult to fully replicate the result of Lee et al. (2020) without being able to easily run the provided source code. In this case, we chose to use measurement model that considers only symptomatic individuals for both phases of the epidemic, as this seemed to visually match the results of Lee et al. (2020) most closely.

B.5.3 Model 3 Replication

Model 3 was developed by a team of researchers at the Laboratory of the Swiss Federal Institute of Technology in Lausanne, hereafter referred to as the Model 3 authors. The code that was originally used to implement Model 3 is archived with the DOI: [10.5281/zenodo.3360723](https://doi.org/10.5281/zenodo.3360723), and also available in the public GitHub repository: `jcblemai/haiti-mass-ocv-campaign`. Because the code was made publicly available, and final model parameters were reported in the supplementary material of Lee et al. (2020), we were able to reproduce Model 3 by directly using the source code. In Fig. B.11, we plot simulations from this model. This figure can be compared to Figure S18 of Lee et al. (2020). We note that slight differences may be accounted for due to variance in the model simulations and the difference in programming language used to produce the figure. Overall, the high standard of reproducibility that was achieved by the Model 3 authors facilitated the ability to readily replicate their model and results.

¹One example of why the code cannot be run that the file loads functions from a non-extant file named `choleraEqs.py` in line 13 rather than `choleraEqsPublic.py`.

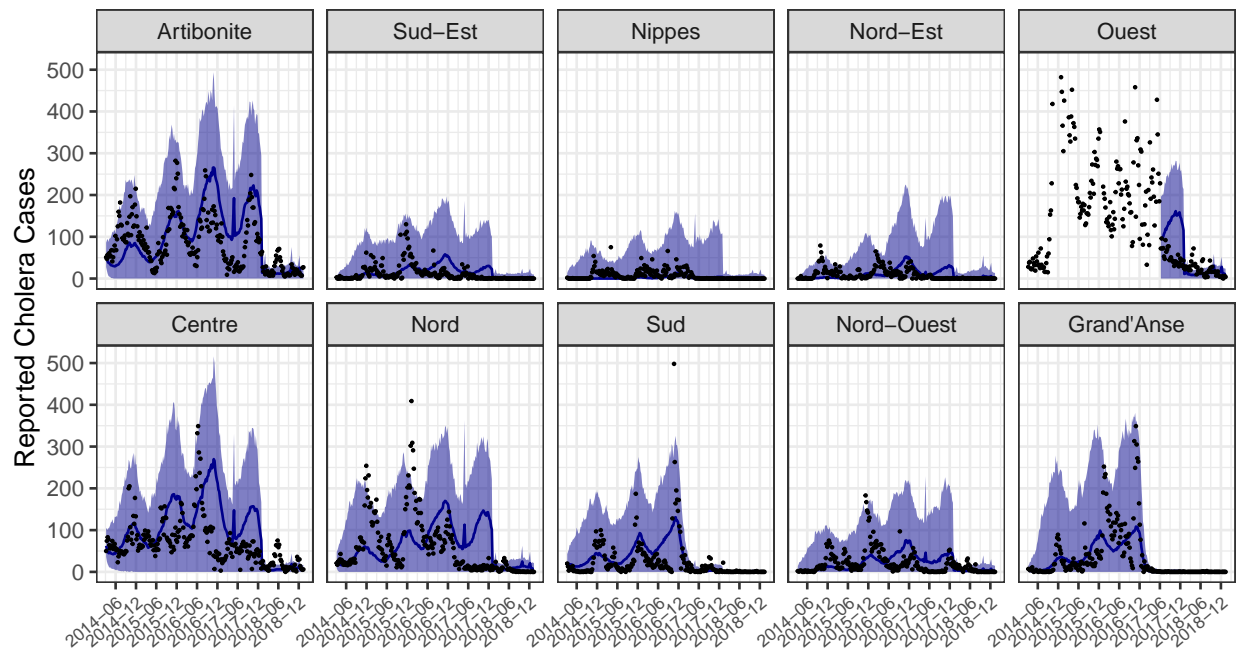


Figure B.11: Simulations from Model 3. Compare to Figure S18 in the supplement of Lee et al. (2020).

B.6 Calibrating Model 3 to observed cases

In this section, we provide more detail on the process that was used to estimate the coefficients of Model 3. In particular, we discuss why we decided to include additional model parameters—those that are associated with the behavior of the system during Hurricane Matthew—that were not considered by Lee et al. (2020). To calibrate this model, we used the iterated block particle filter (IBPF) method of Ionides et al. (2022). Due to the novelty of this algorithm, there exists only a few examples where the IBPF algorithm is used for data analysis Li et al. (2023); Ionides et al. (2022), which is one motivation of the inclusion additional details related to fitting and diagnosing the model fit provided here.

Lee et al. (2020) only estimated model parameters to a simplified version of Model 3 on a subset of the available data, as no method existed at the time of their publication to fit a fully coupled meta-population model to disease incidence data. Building on their results, we fit the fully coupled version of Model 3 to (nearly) all available data, reserving only a few observations to use to calibrate the initial conditions of the model (see the supplement for initialization models for more details). To maximize model likelihoods, we relied on parameter estimates obtained while calculating profile-likelihood confidence intervals, as this calculation requires many replicated IBPF searches. In our preliminary

investigations that were done prior to conducting a profile likelihood search, we found that it was necessary to use multiple searches for the MLE, periodically pruning away less successful searches. To do this, the first collection of searches is performed by obtaining initial values for the parameters by uniformly sampling values from a predefined hypercube. A subsequent refinement search used parameter values corresponding to the largest model likelihoods as starting parameter values. The need for multiple searches does not appear to be uncommon, as a similar approach was used in Ionides et al. (2022). While computationally intensive, profile likelihoods proved to be an effective alternative to maximizing model likelihoods without the need to apply this multistage heuristic.

We use the iterative fitting / pruning technique described above to fit the fully coupled version of Model 3 proposed in Lee et al. (2020). The maximum likelihood we obtained after two rounds of searching was -1.7549×10^4 , which is higher than the benchmark model (-1.7933×10^4). While beating a simple associative benchmark is promising, this does not immediately imply that the model is a good description of the system. Additional investigation of parameters estimates and their corresponding implications on model based conclusions should always be conducted. For meta-population models, it is worth considering how well the model fits the data to each spatial unit. The likelihoods for each department, compared to the corresponding benchmark model, are displayed in Fig. B.12. The figure demonstrates that while the log-likelihood of the fitted model outperforms the auto-regressive negative binomial benchmark model at the aggregate level, Model 3 has lower likelihoods for some departments.

In addition to considering the conditional log-likelihoods of each unit, one can consider conditional log-likelihoods of each observation. When compared to a benchmark, this level of detail can provide useful information about which observations are well described by the model and which are not. In Fig. B.13, we plot the conditional log-likelihoods of Model 3 for each observation. Typically it is most useful to compare the conditional log-likelihoods of the model under consideration to a benchmark, as plotting only conditional log-likelihoods without a comparison may not be helpful. In this case, however, the same insight can be drawn using a figure without a benchmark comparison, so we do not include the benchmark in order to avoid the issue of over-plotting.

Fig. B.13 reveals that the fitted model poorly describes certain features of the data. For example, many departments (in particular Sud) have observations with lower conditional log-likelihoods near October 2016 than at other time points. Further investigation of the model output reveals that the model is struggling to explain the sudden surge in cholera cases that occurred at this time, which coincides with the time that Hurricane Matthew struck Haiti. While the model does include a mechanism to account for increased cholera

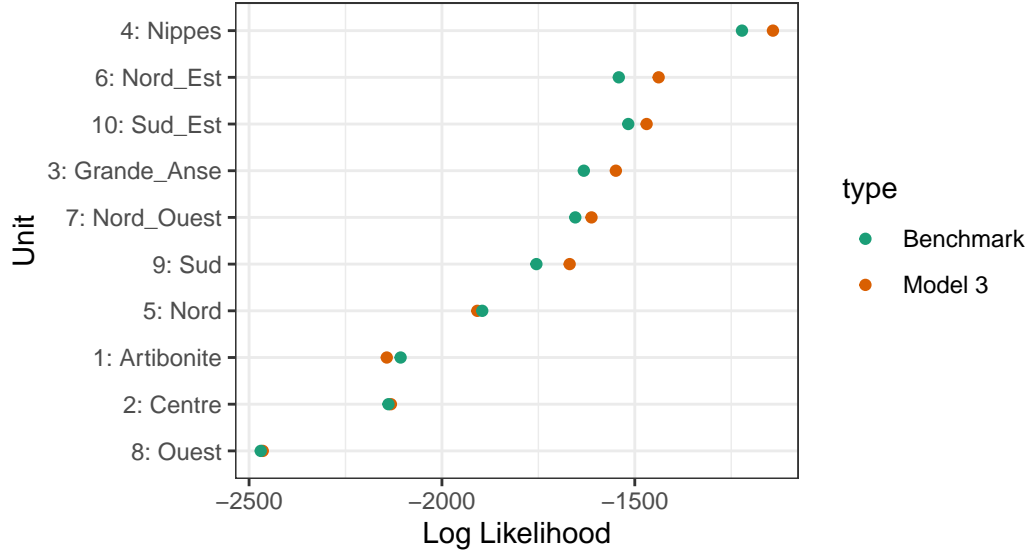


Figure B.12: Log-likelihoods of Model 3 for each department compared to the corresponding benchmark model prior to the inclusion of parameters that account for Hurricane Matthew.

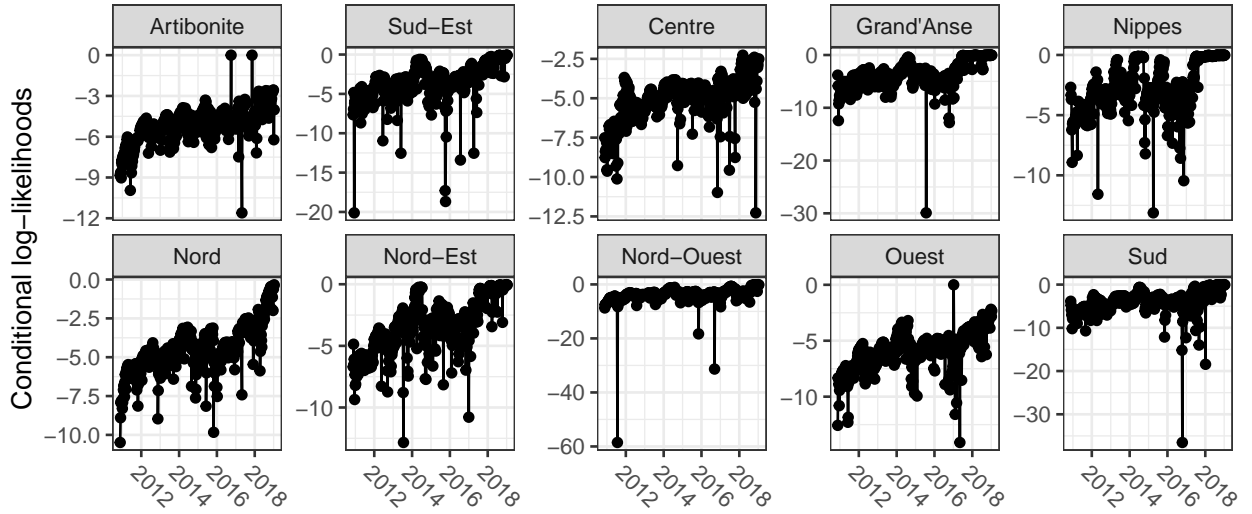


Figure B.13: Conditional log-likelihoods of Model 3 prior to the inclusion of the Hurricane Matthew related parameters.

transmission due to large rainfall events, the mechanism does not appear to be sufficient to capture the damaging effects of the hurricane, which had the greatest impact in the the Sud and Grand'Anse departments Ferreira (2016). This result led us to include parameters β_{Wu}^{hm} and h_u^{hm} in Eq. 23 of the main text, which allows for an increase in the transmission rate between environmental cholera and humans for in Sud and Grand'Anse during and after

the hurricane. The effect of the hurricane on cholera transmission is assumed to have an exponential decay, where the magnitude is determined by β_{Wu}^{hm} and the duration of the effect determined by h_u^{hm} .

We refit Model 3 after introducing these hurricane-related parameters. The resulting model has a log-likelihood value of -1.73329×10^4 . The inclusion of these parameters resulted in an overall increase of 216.4 log-likelihood units. Such a large difference in log-likelihoods is well beyond the threshold of statistical uncertainty determined by Wilks' theorem, suggesting that the data highly favor the inclusion of the additional parameters. The addition of the Hurricane parameters also increases in conditional likelihoods for each observation, particularly around October 2016 (Fig. B.14).

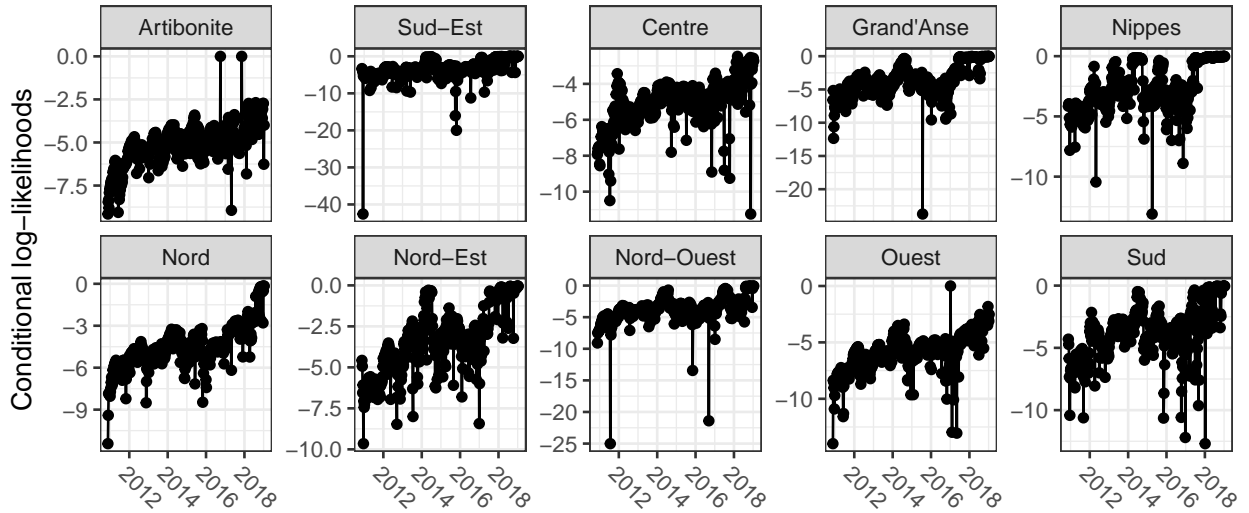


Figure B.14: Conditional log-likelihoods of Model 3 after adding and estimating the parameters related to Hurricane Matthew.

Now that the model with additional parameters has been calibrated to the incidence data, we plot the conditional log-likelihood of each department compared to a benchmark model in Fig. B.15. The difference in log-likelihoods between Model 3 and the benchmark model is smallest in the departments Artibonite, Nord, Ouest and Centre. Each of these departments also exhibited the most sustained cholera transmission, defined by having the fewest number of weeks with no recorded cholera cases. Specifically, these four departments have zero cholera cases recorded in less than 4% of the available data, and all remaining departments—except for Nord-Ouest, which has approximately 9.5% of cases that are zeros and also exhibits the next smallest difference in log-likelihoods—have zero cases recorded in at least 14% of the available weekly data. This result suggests that the quantitative advantage Model 3 has over

its respective benchmark is primarily due to the model’s ability to describe a resurgence of cases after a department records a week of zero cholera cases. This result may be unsurprising in the context of the models that we are comparing. Because the cholera transmission in individual departments likely depends on the national prevalence of cholera and the *Vibrio cholerae* bacteria, our spatially-independent benchmark model that relies exclusively on the previous number of case within any given unit has a difficult time predicting a resurgence of cases.

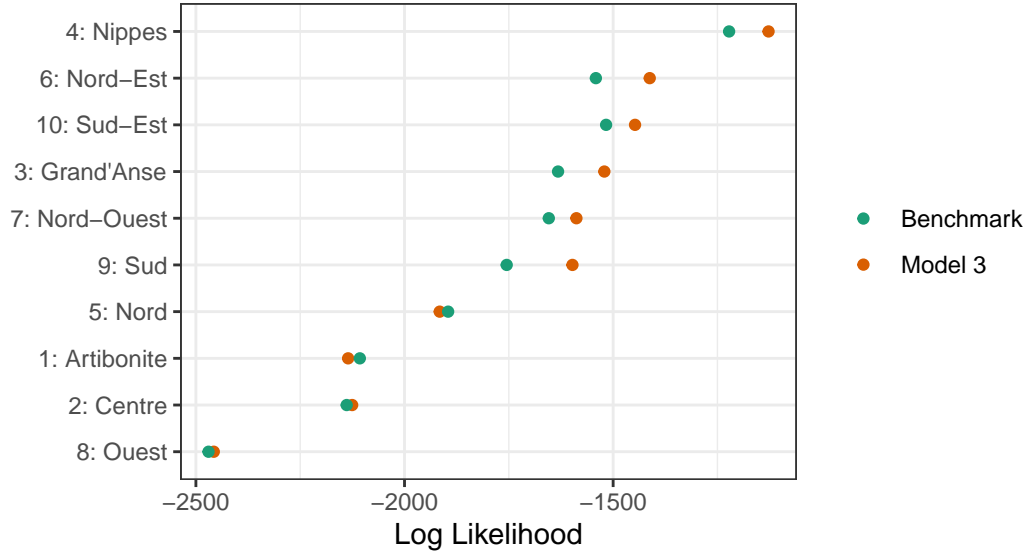


Figure B.15: Log-likelihoods of Model 3 for each department compared to the corresponding benchmark model after adding and estimating parameters related to Hurricane Matthew.

The difference in log-likelihoods between Model 3 and its benchmark model for each individual units suggests that Model 3 has a relatively poor fit for the four units with the most sustained cholera transmission. The simple four parameter benchmark has a higher likelihood than Model 3 for the Artibonite and Nord departments, and also has log-likelihoods only a few units smaller than Model 3 for the departments Ouest and Centre. This is particularly worrisome given that these four departments account for more than 77% of the total number of reported cholera cases.

B.6.1 Examining the Hidden States of the Calibrated Model

For mechanistic models, beating a suitable statistical benchmark does not alone guarantee that the model provides an accurate description of a dynamic process. Indeed, a good statistical fit does not require the model to assert a causal explanation. For example, reconstructed latent variables should make sense in the context of alternative measurements of these variables

Grad et al. (2012). We demonstrate this principle by examining the latent state of the calibrated model. In particular, we examine the compartment of susceptible individuals under various scenarios. This analysis can also provide insight into why the calibrated model fails to outperform the benchmark model on the four departments with the most sustained cholera transmission.

Recall that the filtering distribution for the calibrated version of Model 3 at time t_k is defined as the distribution of the latent state at time t_k given the data from times $t_1 : t_k$, i.e. $f_{\mathbf{X}_k | \mathbf{Y}_{1:k}}^{(3)}(\mathbf{x}_k | \mathbf{y}_{1:k}^*; \hat{\theta})$. In general, one may expect simulations from the filtering distribution of a model with a good statistical fit to result in hidden states that are highly consistent with the observed data because the filtering distribution is conditioned on the observed data. Fig. B.16 shows the percentage of individuals that are in the susceptible compartment from various simulations of the model: simulations from Model 3 under initial conditions are displayed in red; simulations from the filtering distribution of model are displayed in blue. This figure shows that simulations from initial conditions tends to result in a much more rapid depletion of susceptible individuals at the start of the epidemic than simulations from the filtering distribution, suggesting the calibrated model has a propensity to predict larger outbreaks than what is typically seen in the data. This result demonstrates that the calibrated model favors a more rapid growth in cholera cases than what is typically seen in the observed data, providing a possible explanation as to why the model fails to beat the simple benchmark for each spatial unit. This results hints at the possibility of model misspecification, and warrants a degree of caution in interpreting the model's outputs.

B.7 Forecasting with parameter uncertainty

Let $f_{Y_{1:N}}(y_{1:N} | \theta)$ denote the pdf of the model under consideration, where θ is a parameter vector that indexes the model. Furthermore, denote the observed data as $y_{1:N}^*$. Because the uncertainty in just a single parameter can lead to drastically different forecasts Saltelli et al. (2020), parameter uncertainty should be considered when obtaining model forecasts when the goal is to influencing policy. In a Bayesian modelling paradigm, the most natural way to account for parameter uncertainty in model forecasts is to suppose that θ comes from a distribution f_{Θ} , and then to obtain J forecasts from the model where each forecast is obtained using parameters drawn from the posterior distribution $\theta_{1:J} | Y_{1:N} = y_{1:N}^* \sim f_{\Theta}(\theta | Y_{1:N} = y_{1:N}^*)$.

When frequentist methods are used, however, there does not exist a posterior distribution from which one could sample. A common approach could be to obtain a weighted average of the simulations from various models Hoeting et al. (1999), but this can be problematic

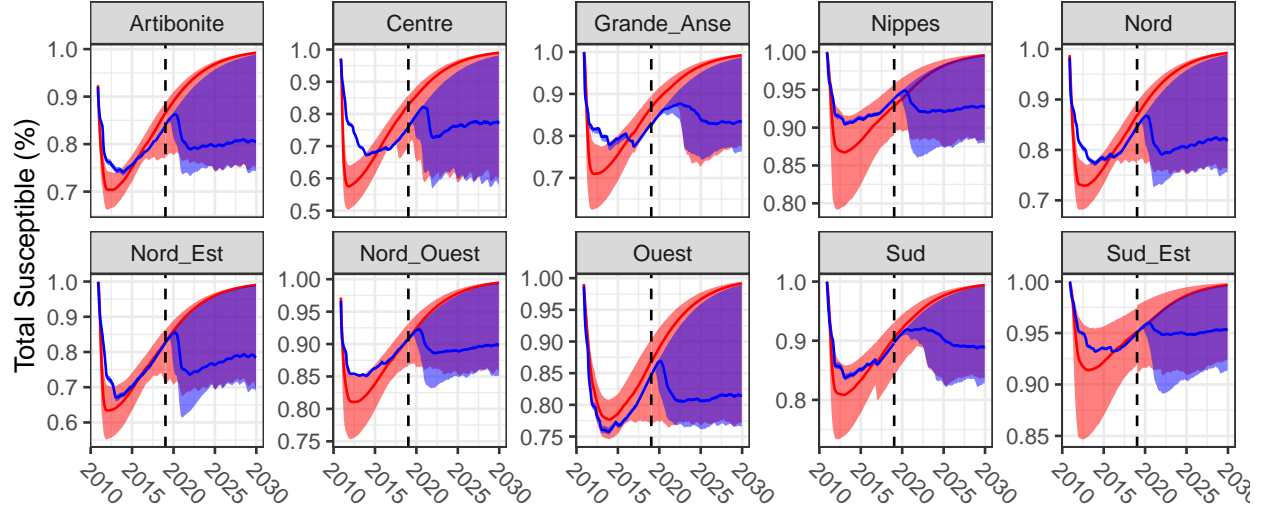


Figure B.16: Percentage of individuals that are in the susceptible compartment. Simulations from Model 3 under initial conditions are displayed in red; simulations from the filtering distribution of model are displayed in blue. The dashed line represents the end of the observed data.

when forecasts from each model are very different from each other Grueber et al. (2011). A similar approach that has been taken King et al. (2015) is to obtain model forecasts using multiple sets of parameter values and then sample from the resulting forecasts using weights proportional to the corresponding likelihoods of the parameter values. This approach could be considered as empirical Bayes, as it is equivalent to using a discrete uniform prior where the set of values in the prior distribution is determined by a stochastic routine applied to the observed data, as discussed below.

For each $k \in 1 : K$, let θ_k be a unique set of model parameters. Letting Θ denote a random vector of model parameters, we endow Θ with a discrete uniform distribution on the set $\{\theta_1, \theta_2, \dots, \theta_K\}$, such that $P(\Theta = \theta_k) = \frac{1}{K}$ for all values $k \in 1 : K$. Using this as a prior distribution, the posterior distribution of $\Theta | Y_{1:N} = y_{1:N}^*$ can be expressed as: $P(\Theta = \theta_k | Y_{1:N} = y_{1:N}^*) = \frac{f_{Y_{1:N}}(y_{1:N}^* | \theta_k)}{\sum_{l=1}^K f_{Y_{1:N}}(y_{1:N}^* | \theta_l)}$. In a standard empirical Bayes analysis, the values $\theta_1, \dots, \theta_K$ of the prior distribution would be chosen using the observed data, resulting in a posterior distribution that weighs the prior parameter vectors proportional to their corresponding likelihoods. We choose θ_k to be the output of a stochastic routine applied to the observed data by setting θ_k to be the output of an iterated filtering algorithm. In practice, because the likelihood maximization routines of iterated filtering methods are stochastic, it is common to run the iterated filtering method multiple times (K) for each model in order to obtain a maximum likelihood estimate for model parameters. This results in a natural set of

parameters near the MLE that could be used as the discrete prior distribution. Alternatively, the set $\{\theta_1, \theta_2, \dots, \theta_K\}$ could be determined by first obtaining marginal confidence intervals for each element of the parameter vector Θ , and then creating a hypercube using the combination of marginal confidence intervals. The set $\{\theta_1, \theta_2, \dots, \theta_K\}$ is then obtained by sampling uniformly K values from the resulting hypercube, as was done by King et al. (2015).

APPENDIX C

Example Appendix 01

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

C.1 Sample appendix section

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

C.1.1 Sample appendix subsection

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in

sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

C.2 Another sample appendix section

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

APPENDIX D

Example Appendix 02

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

BIBLIOGRAPHY

- Abrams, J. Y., J. R. Copeland, R. V. Tauxe, K. A. Date, E. D. Belay, R. K. Mody, and E. D. Mintz (2013). Real-time modelling used for outbreak management during a cholera epidemic, haiti, 2010–2011. *Epidemiology and Infection* 141(6), 1276–1285.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19(6), 716–723.
- Akman, O. and E. Schaefer (2015). An evolutionary computing approach for parameter estimation investigation of a model for cholera. *Journal of biological dynamics* 9(1), 147–158.
- Andrews, J. R. and S. Basu (2011). Transmission dynamics and control of cholera in Haiti: an epidemic model. *The Lancet* 377(9773), 1248–1255.
- Andrieu, C., A. Doucet, and R. Holenstein (2010). Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* 72, 269–342.
- Arulampalam, M. S., S. Maskell, N. Gordon, and T. Clapp (2002). A tutorial on particle filters for online nonlinear, non-Gaussian Bayesian tracking. *IEEE Transactions on Signal Processing* 50, 174 – 188.
- Asfaw, K., J. Park, A. A. King, and E. L. Ionides (2024). spatpomp: An R package for spatiotemporal partially observed Markov process models. *Journal of Open Source Software* 9(104), 7008.
- Azman, A. S., F. J. Luquero, I. Ciglenecki, R. F. Grais, D. A. Sack, and J. Lessler (2015). The impact of a one-dose versus two-dose oral cholera vaccine regimen in outbreak settings: A modeling study. *PLoS Medicine* 12(8), e1001867.
- Azman, A. S., F. J. Luquero, A. Rodrigues, P. P. Palma, R. F. Grais, C. N. Banga, B. T. Grenfell, and J. Lessler (2012). Urban cholera transmission hotspots and their implications for reactive vaccination: evidence from Bissau city, Guinea bissau. *PLoS neglected tropical diseases* 6(11), e1901.
- Baker, R. E., J.-M. Pena, J. Jayamohan, and A. Jérusalem (2018). Mechanistic models versus machine learning, a fight worth fighting for the biological community? *Biology letters* 14(5), 20170660.

- Barzilay, E. J., N. Schaad, R. Magloire, K. S. Mung, J. Boncy, G. A. Dahourou, E. D. Mintz, M. W. Steenland, J. F. Vertefeuille, and J. W. Tappero (2013). Cholera surveillance during the Haiti epidemic—the first 2 years. *New England Journal of Medicine* 368(7), 599–609. Includes case definition to 2013: The NCSS used a modified World Health Organization case definition for cholera that included acute watery diarrhea, with or without vomiting, in persons of all ages residing in an area in which at least one case of *Vibrio cholerae* O1 infection had been confirmed by culture.
- Behrend, M. R., M.-G. Basáñez, J. I. D. Hamley, T. C. Porco, W. A. Stolk, M. Walker, S. J. de Vlas, and for the NTD Modelling Consortium (2020, 04). Modelling for policy: The five principles of the neglected tropical diseases modelling consortium. *PLoS Neglected Tropical Diseases* 14(4), 1–17.
- Bengtsson, L., J. Gaudart, X. Lu, S. Moore, E. Wetter, K. Sallah, S. Rebaudet, and R. Piarroux (2015). Using mobile phone data to predict the spatial spread of cholera. *Scientific reports* 5(1), 8923.
- Bhadra, A., E. L. Ionides, K. Laneri, M. Pascual, M. Bouma, and R. 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, 440–451.
- Botelho, C., J. D. Kong, M. A. Lucien, Z. Shuai, and H. Wang (2021). A mathematical model for *Vibrio*-phage interactions. *Mathematical Biosciences and Engineering* 18(3).
- Box, G. and G. Jenkins (1970). *Time Series Analysis: Forecasting and Control*. San Francisco, Holden-Day.
- Box, G. E. (1979). Robustness in the strategy of scientific model building. In *Robustness in Statistics*, pp. 201–236. Elsevier.
- Brauer, F. (2017). Mathematical epidemiology: Past, present, and future. *Infectious Disease Modelling* 2(2), 113–127.
- Bretó, C., D. He, E. L. Ionides, and A. A. King (2009). Time series analysis via mechanistic models. *Annals of Applied Statistics* 3, 319–348.
- Bretó, C. and E. L. Ionides (2011). Compound Markov counting processes and their applications to modeling infinitesimally over-dispersed systems. *Stochastic Processes and their Applications* 121, 2571–2591.
- Brockwell, P. J. and R. A. Davis (1991). *Time Series: Theory and Methods*. Springer Series in Statistics. New York, NY: Springer.
- Capone, F., V. De Cataldis, and R. De Luca (2015). Influence of diffusion on the stability of equilibria in a reaction–diffusion system modeling cholera dynamic. *Journal of mathematical biology* 71, 1107–1131.

- Chakhchoukh, Y. (2010). A new robust estimation method for arma models. *IEEE Transactions on Signal Processing* 58(7), 3512–3522.
- Chao, D. L., M. E. Halloran, and I. M. Longini Jr (2011). Vaccination strategies for epidemic cholera in haiti with implications for the developing world. *Proceedings of the National Academy of Sciences* 108(17), 7081–7085.
- Charles, M., G. G. Delva, J. Boutin, K. Severe, M. Peck, M. M. Mabou, P. F. Wright, and J. W. Pape (2014). Importance of cholera and other etiologies of acute diarrhea in post-earthquake port-au-prince, haiti. *The American journal of tropical medicine and hygiene* 90(3), 511.
- Chib, S. and E. Greenberg (1994). Bayes inference in regression models with arma (p, q) errors. *Journal of Econometrics* 64(1), 183–206.
- Chon, K. and R. Cohen (1997). Linear and nonlinear arma model parameter estimation using an artificial neural network. *IEEE Transactions on Biomedical Engineering* 44(3), 168–174.
- Collins, O. and K. Govinder (2014). Incorporating heterogeneity into the transmission dynamics of a waterborne disease model. *Journal of Theoretical Biology* 356, 133–143.
- Collins, O. C. and K. J. Duffy (2021). Mathematical analyses on the effects of control measures for a waterborne disease model with socioeconomic conditions. *Journal of Computational Biology* 28(1), 19–32.
- Cuneo, C. N., E. Dansereau, A. R. Habib, M. Davies, S. Ware, and K. Kornetsky (2017). Treating childhood malnutrition in rural haiti: Program outcomes and obstacles. *Annals of global health* 83(2), 300–310.
- Dadlani, A., R. O. Afolabi, H. Jung, K. Sohraby, and K. Kim (2020). Deterministic models in epidemiology: From modeling to implementation. *arXiv:2004.04675*.
- Dahabreh, I. J., J. A. Chan, A. Earley, D. Moorthy, E. E. Avendano, T. A. Trikalinos, E. M. Balk, and J. B. Wong (2017). Modeling and simulation in the context of health technology assessment: Review of existing guidance, future research needs, and validity assessment.
- Date, K. A., A. Vicari, T. B. Hyde, E. Mintz, M. C. Danovaro-Holliday, A. Henry, J. W. Tappero, T. H. Roels, J. Abrams, B. T. Burkholder, et al. (2011). Considerations for oral cholera vaccine use during outbreak after earthquake in haiti, 2010- 2011. *Emerging infectious diseases* 17(11), 2105.
- Donnelly, C. A., I. Boyd, P. Campbell, C. Craig, P. Vallance, M. Walport, C. J. Whitty, E. Woods, and C. Wormald (2018). Four principles to make evidence synthesis more useful for policy.
- Donoho, D. (2017). 50 years of data science. *Journal of Computational and Graphical Statistics* 26(4), 745–766.

- Doucet, A., N. De Freitas, N. J. Gordon, et al. (2001). *Sequential Monte Carlo methods in practice*, Volume 1. Springer.
- Durbin, J. and S. J. Koopman (2012). *Time Series Analysis by State Space Methods*, Volume 38. OUP Oxford.
- Eddy, S. R. (2004). What is a hidden markov model? *Nature biotechnology* 22(10), 1315–1316.
- Egger, M., L. Johnson, C. Althaus, A. Schöni, G. Salanti, N. Low, and S. L. Norris (2017). Developing who guidelines: time to formally include evidence from mathematical modelling studies. *F1000Research* 6.
- Eisenberg, M. C., G. Kujbida, A. R. Tuite, D. N. Fisman, and J. H. Tien (2013). Examining rainfall and cholera dynamics in haiti using statistical and dynamic modeling approaches. *Epidemics* 5(4), 197–207.
- Evensen, G. (2009). The ensemble Kalman filter for combined state and parameter estimation. *IEEE Transactions on Control Systems* 29, 83–104.
- Evensen, G., F. C. Vossepoel, and P. J. van Leeuwen (2022). *Data assimilation fundamentals: A unified formulation of the state and parameter estimation problem*. Springer Nature.
- Ferreira, S. (2016). Cholera threatens Haiti after Hurricane Matthew. *BMJ* 355, i5516.
- Fitzgibbon, W. E., J. J. Morgan, G. F. Webb, and Y. Wu (2020). Modelling the aqueous transport of an infectious pathogen in regional communities: application to the cholera outbreak in Haiti. *Journal of the Royal Society Interface* 17(169), 20200429.
- Fletcher, R. (2000). *Practical Methods of Optimization*. John Wiley & Sons.
- Francois, J. (2020). Cholera remains a public health threat in haiti. *The Lancet Global Health* 8(8), e984.
- Fung, I. C.-H., D. L. Fitter, R. H. Borse, M. I. Meltzer, and J. W. Tappero (2013). Modeling the effect of water, sanitation, and hygiene and oral cholera vaccine implementation in haiti. *The American Journal of Tropical Medicine and Hygiene* 89(4), 633.
- Ganusov, V. V. (2016). Strong inference in mathematical modeling: a method for robust science in the twenty-first century. *Frontiers in Microbiology* 7, 1131.
- Gardner, G., A. C. Harvey, and G. D. A. Phillips (1980). Algorithm AS 154: An algorithm for exact maximum likelihood estimation of autoregressive-moving average models by means of Kalman filtering. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 29(3), 311–322.
- Gatto, M., L. Mari, E. Bertuzzo, R. Casagrandi, L. Righetto, I. Rodriguez-Iturbe, and A. Rinaldo (2012). Generalized reproduction numbers and the prediction of patterns in waterborne disease. *Proceedings of the National Academy of Sciences* 109(48), 19703–19708.

- Gentleman, R. and D. Temple Lang (2007). Statistical analyses and reproducible research. *Journal of Computational and Graphical Statistics* 16(1), 1–23.
- Giordano, G., F. Blanchini, R. Bruno, P. Colaneri, A. Di Filippo, A. Di Matteo, and M. Colaneri (2020). Modelling the covid-19 epidemic and implementation of population-wide interventions in italy. *Nature Medicine* 26(6), 855–860.
- Glennie, R., T. Adam, V. Leos-Barajas, T. Michelot, T. Photopoulou, and B. T. McClintock (2023). Hidden Markov models: Pitfalls and opportunities in ecology. *Methods in Ecology and Evolution* 14(1), 43–56.
- Grad, Y. H., J. C. Miller, and M. Lipsitch (2012). Cholera modeling: Challenges to quantitative analysis and predicting the impact of interventions. *Epidemiology* 23(4), 523.
- Green, K. C. and J. S. Armstrong (2015). Simple versus complex forecasting: The evidence. *Journal of Business Research* 68(8), 1678–1685. Special Issue on Simple Versus Complex Forecasting.
- Grueber, C. E., S. Nakagawa, R. J. Laws, and I. G. Jamieson (2011). Multimodel inference in ecology and evolution: challenges and solutions. *Journal of Evolutionary Biology* 24(4), 699–711.
- Hart, W. E. (1998). Sequential stopping rules for random optimization methods with applications to multistart local search. *SIAM Journal on Optimization* 9(1), 270–290.
- He, D., E. L. Ionides, and A. A. King (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.
- Henrys, J. H., G. Lerebours, M. A. Achille, K. Moise, and C. Raccurt (2020). Cholera in haiti. *The Lancet Global Health* 8(12), e1469.
- Hoeting, J. A., D. Madigan, A. E. Raftery, and C. T. Volinsky (1999). Bayesian model averaging: a tutorial (with comments by M. Clyde, David Draper and E. I. George, and a rejoinder by the authors. *Statistical Science* 14(4), 382 – 417.
- Hulland, E., S. Subaiya, K. Pierre, N. Barthelemy, J. S. Pierre, A. Dismer, S. Juin, D. Fitter, and J. Brunkard (2019). Increase in reported cholera cases in haiti following hurricane matthew: an interrupted time series model. *The American journal of tropical medicine and hygiene* 100(2), 368.
- Hyndman, R. J. and Y. Khandakar (2008). Automatic time series forecasting: The forecast package for R. *Journal of Statistical Software* 26(3), 1–22.
- Ioannidis, J. P., S. Cripps, and M. A. Tanner (2020). Forecasting for COVID-19 has failed. *International Journal of Forecasting*.
- Ionides, E. L., K. Asfaw, J. Park, and A. A. King (2021). Bagged filters for partially observed interacting systems. *Journal of the American Statistical Association* pre-published online.

- Ionides, E. L., C. Breto, J. Park, R. A. Smith, and A. A. King (2017). Monte Carlo profile confidence intervals for dynamic systems. *Journal of the Royal Society Interface* 14, 1–10.
- Ionides, E. L., D. Nguyen, Y. Atchadé, S. Stoev, and A. A. King (2015). Inference for dynamic and latent variable models via iterated, perturbed Bayes maps. *Proceedings of the National Academy of Sciences of the USA* 112(3), 719—724.
- Ionides, E. L., N. Ning, and J. Wheeler (2022). An iterated block particle filter for inference on coupled dynamic systems with shared and unit-specific parameters. *Statistica Sinica*, pre-published online.
- Ivers, L. C. (2017). Eliminating cholera transmission in haiti. *New England Journal of Medicine* 376(2), 101–103.
- Ivers, L. C., I. J. Hilaire, J. E. Teng, C. P. Almazor, J. G. Jerome, R. Ternier, J. Boncy, J. Buteau, M. B. Murray, J. B. Harris, et al. (2015). Effectiveness of reactive oral cholera vaccination in rural haiti: a case-control study and bias-indicator analysis. *The Lancet Global Health* 3(3), e162–e168.
- Kalman, R. E. (1960, 03). A new approach to linear filtering and prediction problems. *Journal of Basic Engineering* 82(1), 35–45.
- Kelly Jr, M. R., J. H. Tien, M. C. Eisenberg, and S. Lenhart (2016). The impact of spatial arrangements on epidemic disease dynamics and intervention strategies. *Journal of biological dynamics* 10(1), 222–249.
- Kermack, W. O. and A. G. McKendrick (1927). A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London, Series A* 115(772), 700–721.
- King, A. A., M. Domenech de Cellès, F. M. Magpantay, and P. Rohani (2015). Avoidable errors in the modelling of outbreaks of emerging pathogens, with special reference to ebola. *Proceedings of the Royal Society B: Biological Sciences* 282(1806), 20150347.
- King, A. A., D. Nguyen, and E. L. Ionides (2016). Statistical inference for partially observed Markov processes via the R package pomp. *Journal of Statistical Software* 69, 1–43.
- King, A. A. and T. Rowan (2020). **subplex**: Unconstrained optimization using the subplex algorithm. R package, available at <https://cran.r-project.org/web/packages/subplex>.
- Kirpich, A., T. A. Weppelmann, Y. Yang, A. Ali, J. G. Morris, Jr., and I. M. Longini (2015, 10). Cholera transmission in ouest department of haiti: Dynamic modeling and the future of the epidemic. *PLOS Neglected Tropical Diseases* 9(10), 1–12.
- Kirpich, A., T. A. Weppelmann, Y. Yang, J. G. Morris Jr, and I. M. Longini Jr (2017). Controlling cholera in the ouest department of haiti using oral vaccines. *PLOS Neglected Tropical Diseases* 11(4), e0005482.

- Kühn, J., F. Finger, E. Bertuzzo, S. Borgeaud, M. Gatto, A. Rinaldo, and M. Blokesch (2014). Glucose-but not rice-based oral rehydration therapy enhances the production of virulence determinants in the human pathogen *Vibrio cholerae*. *PLoS neglected tropical diseases* 8(12), e3347.
- Kunkel, A., J. A. Lewnard, V. E. Pitzer, and T. Cohen (2017). Antimicrobial resistance risks of cholera prophylaxis for United Nations peacekeepers. *Antimicrobial agents and chemotherapy* 61(8), 10–1128.
- Laneri, K., A. Bhadra, E. L. Ionides, M. Bouma, R. Yadav, R. Dhiman, and M. Pascual (2010). Forcing versus feedback: Epidemic malaria and monsoon rains in NW India. *PLoS Computational Biology* 6, e1000898.
- Lang, M., B. Bischl, and D. Surmann (2017, feb). batchtools: Tools for r to work on batch systems. *The Journal of Open Source Software* (10).
- Lau, M. S., A. Becker, W. Madden, L. A. Waller, C. J. E. Metcalf, and B. T. Grenfell (2022). Comparing and linking machine learning and semi-mechanistic models for the predictability of endemic measles dynamics. *PLoS computational biology* 18(9), e1010251.
- Le Cam, L. (1990). Maximum likelihood: An introduction. *International Statistical Review / Revue Internationale de Statistique* 58(2), 153–171.
- Lee, E. C., D. L. Chao, J. C. Lemaitre, L. Matrajt, D. Pasetto, J. Perez-Saez, F. Finger, A. Rinaldo, J. D. Sugimoto, M. E. Halloran, I. M. Longini, R. Ternier, K. Vissieres, A. S. Azman, J. Lessler, and L. C. Ivers (2020). Achieving coordinated national immunity and cholera elimination in Haiti through vaccination: A modelling study. *The Lancet Global Health* 8(8), e1081–e1089.
- Lee, E. C., R. Ternier, J. Lessler, A. S. Azman, and L. C. Ivers (2020). Cholera in haiti—authors’ reply. *The Lancet Global Health* 8(12), e1470–e1471.
- Leung, T., J. Eaton, and L. Matrajt (2022). Optimizing one-dose and two-dose cholera vaccine allocation in outbreak settings: A modeling study. *PLOS Neglected Tropical Diseases* 16(4), e0010358.
- Lewis, A. S. L., C. R. Rollinson, A. J. Allyn, J. Ashander, S. Brodie, C. B. Brookson, E. Collins, M. C. Dietze, A. S. Gallinat, N. Juvigny-Khenafou, G. Koren, D. J. McGlinn, H. Moustahfid, J. A. Peters, N. R. Record, C. J. Robbins, J. Tonkin, and G. M. Wardle (2022). The power of forecasts to advance ecological theory. *Methods in Ecology and Evolution*, pre-published online.
- Lewnard, J. A., M. Antillón, G. Gonsalves, A. M. Miller, A. I. Ko, and V. E. Pitzer (2016). Strategies to prevent cholera introduction during international personnel deployments: a computational modeling analysis based on the 2010 haiti outbreak. *PLoS medicine* 13(1), e1001947.
- Li, J., E. L. Ionides, A. A. King, M. Pascual, and N. Ning (2023). Machine learning for mechanistic models of metapopulation dynamics. *arxiv:2311.06702*.

- Li, L., B. A. Brumback, T. A. Weppelmann, J. G. Morris Jr, and A. Ali (2016). Adjusting for unmeasured confounding due to either of two crossed factors with a logistic regression model. *Statistics in Medicine* 35(18), 3179–3188.
- Lii, K.-S. (1990). Identification and estimation of non-gaussian arma processes. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 38(7), 1266–1276.
- Lin, J., R. Xu, and X. Tian (2019). Transmission dynamics of cholera with hyperinfectious and hypoinfectious vibrios: mathematical modelling and control strategies. *Mathematical Biosciences and Engineering* 16(5), 4339–4358.
- Liu, J. S. (2001). *Monte Carlo Strategies in Scientific Computing*. New York: Springer.
- Lofgren, E. T., M. E. Halloran, C. M. Rivers, J. M. Drake, T. C. Porco, B. Lewis, W. Yang, A. Vespignani, J. Shaman, J. N. Eisenberg, et al. (2014). Mathematical models: A key tool for outbreak response. *Proceedings of the National Academy of Sciences* 111(51), 18095–18096.
- Lucas, R. E. et al. (1976). Econometric policy evaluation: A critique. In *Carnegie-Rochester Conference Series on Public Policy*, Volume 1, pp. 19–46.
- Mari, L., E. Bertuzzo, F. Finger, R. Casagrandi, M. Gatto, and A. Rinaldo (2015). On the predictive ability of mechanistic models for the Haitian cholera epidemic. *Journal of the Royal Society Interface* 12(104), 20140840.
- Matias, W. R., J. E. Teng, I. J. Hilaire, J. B. Harris, M. F. Franke, and L. C. Ivers (2017). Household and individual risk factors for cholera among cholera vaccine recipients in rural haiti. *The American Journal of Tropical Medicine and Hygiene* 97(2), 436.
- Mavian, C., T. K. Paisie, M. T. Alam, C. Browne, V. M. Beau De Rochars, S. Nembrini, M. N. Cash, E. J. Nelson, T. Azarian, A. Ali, et al. (2020). Toxigenic *Vibrio cholerae* evolution and establishment of reservoirs in aquatic ecosystems. *Proceedings of the National Academy of Sciences* 117(14), 7897–7904.
- May, R. M. (2004). Uses and abuses of mathematics in biology. *science* 303(5659), 790–793.
- Mayo, D. G. (2018). *Statistical Inference as Severe Testing*. Cambridge: Cambridge University Press.
- McCabe, R. and C. A. Donnelly (2021). Disease transmission and control modelling at the science–policy interface. *Interface Focus* 11(6), 20210013.
- Michel, E., J. Gaudart, S. Beaulieu, G. Bulit, M. Piarroux, J. Boncy, P. Dely, R. Piarroux, and S. Rebaudet (2019). Estimating effectiveness of case-area targeted response interventions against cholera in Haiti. *Elife* 8, e50243.
- Moise, K., A. M. Achille, D. Batumbo, B. Bourdeau, S. Rebaudet, G. Lerebours, J. H. Henrys, and C. Raccurt (2020). Impact of patron saint festivities on cholera in three communes in haiti. *BMC Public Health* 20, 1–7.

- Monahan, J. F. (1983). Fully bayesian analysis of arma time series models. *Journal of Econometrics* 21(3), 307–331.
- Mukandavire, Z. and J. G. Morris Jr (2015). Modeling the epidemiology of cholera to prevent disease transmission in developing countries. *Microbiology spectrum* 3(3), 10–1128.
- Mukandavire, Z., D. L. Smith, and J. G. Morris Jr (2013). Cholera in haiti: reproductive numbers and vaccination coverage estimates. *Scientific reports* 3(1), 997.
- Ndii, M. Z. and A. K. Supriatna (2017). Stochastic mathematical models in epidemiology. *Information* 20, 6185–6196.
- Newman, K., R. King, V. Elvira, P. de Valpine, R. S. McCrea, and B. J. Morgan (2023). State-space models for ecological time-series data: Practical model-fitting. *Methods in Ecology and Evolution* 14(1), 26–42.
- Ning, N. and E. L. Ionides (2023). Iterated block particle filter for high-dimensional parameter learning: Beating the curse of dimensionality. *Journal of Machine Learning Research* 24, 1–76.
- NOAA (2016). Monthly average master gauge water levels (1860-present): Lake michigan-huron. Accessed: Jan 24, 2016.
- Pan American Health Organization (2023). Cholera epidemic in hispaniola 2023 - situation report 19.
- Park, J. and E. L. Ionides (2020). Inference on high-dimensional implicit dynamic models using a guided intermediate resampling filter. *Statistics & Computing* 30, 1497–1522.
- Pasetto, D., F. Finger, A. Camacho, F. Grandesso, S. Cohuet, J. C. Lemaitre, A. S. Azman, F. J. Luquero, E. Bertuzzo, and A. Rinaldo (2018, 05). Near real-time forecasting for cholera decision making in haiti after hurricane matthew. *PLOS Computational Biology* 14(5), 1–22.
- Pawitan, Y. (2001). *In all likelihood: statistical modelling and inference using likelihood*. Oxford University Press.
- Pezzoli, L. (2020). Global oral cholera vaccine use, 2013–2018. *Vaccine* 38, A132–A140.
- Peñaloza Ramos, M. C., P. Barton, S. Jowett, and A. J. Sutton (2015). A systematic review of research guidelines in decision-analytic modeling. *Value in Health* 18(4), 512–529.
- Piarroux, R., R. Barraï, B. Faucher, R. Haus, M. Piarroux, J. Gaudart, R. Magloire, and D. Raoult (2011). Understanding the cholera epidemic, haiti. *Emerging infectious diseases* 17(7), 1161.
- Pons-Salort, M. and N. C. Grassly (2018). Serotype-specific immunity explains the incidence of diseases caused by human enteroviruses. *Science* 361(6404), 800–803.

- Prosperi, M., Y. Guo, M. Sperrin, J. S. Koopman, J. S. Min, X. He, S. Rich, M. Wang, I. E. Buchan, and J. Bian (2020). Causal inference and counterfactual prediction in machine learning for actionable healthcare. *Nature Machine Intelligence* 2(7), 369–375.
- R Core Team (2022). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Raila, E. M. and D. O. Anderson (2017). Healthcare waste management during disasters and its effects on climate change: Lessons from 2010 earthquake and cholera tragedies in haiti. *Waste Management & Research* 35(3), 236–245.
- Rebaudet, S., G. Bulit, J. Gaudart, E. Michel, P. Gazin, C. Evers, S. Beaulieu, A. A. Abedi, L. Osei, R. Barraïs, et al. (2019). The case-area targeted rapid response strategy to control cholera in haiti: A four-year implementation study. *PLoS Neglected Tropical Diseases* 13(4), e0007263.
- Rebaudet, S., P. Dély, J. Boncy, J. H. Henrys, and R. Piarroux (2021). Toward cholera elimination, haiti. *Emerging Infectious Diseases* 27(11), 2932.
- Rebaudet, S., J. Gaudart, and R. Piarroux (2020). Cholera in haiti. *The Lancet Global Health* 8(12), e1468.
- Rebeschini, P. and R. van Handel (2015). Can local particle filters beat the curse of dimensionality? *The Annals of Applied Probability* 25(5), 2809–2866.
- Richterman, A., A. S. Azman, G. Constant, and L. C. Ivers (2019). The inverse relationship between national food security and annual cholera incidence: a 30-country analysis. *BMJ Global Health* 4(5).
- Richterman, A., M. F. Franke, G. Constant, G. Jerome, R. Ternier, and L. C. Ivers (2019). Food insecurity and self-reported cholera in haitian households: An analysis of the 2012 demographic and health survey. *PLoS neglected tropical diseases* 13(1), e0007134.
- Rinaldo, A., E. Bertuzzo, L. Mari, L. Righetto, M. Blokesch, M. Gatto, R. Casagrandi, M. Murray, S. M. Vesenbeckh, and I. Rodriguez-Iturbe (2012). Reassessment of the 2010–2011 Haiti cholera outbreak and rainfall-driven multiseason projections. *Proceedings of the National Academy of Sciences* 109(17), 6602–6607.
- Ripley, B. D. (2002, June). Time series in R 1.5.0. *The Newsletter of the R Project Volume* 2, 2.
- Romero-Severson, E., E. Volz, J. Koopman, T. Leitner, and E. Ionides (2015). Dynamic variation in sexual contact rates in a cohort of HIV-negative gay men. *American Journal of Epidemiology* 182, 255–262.
- Rubin, D. H. F., F. G. Zingl, D. R. Leitner, R. Ternier, V. Compere, S. Marseille, D. Slater, J. B. Harris, F. Chowdhury, F. Qadri, J. Boncy, L. C. Ivers, and M. K. Waldor (2022). Reemergence of cholera in haiti. *New England Journal of Medicine*.

- Sallah, K., R. Giorgi, L. Bengtsson, X. Lu, E. Wetter, P. Adrien, S. Rebaudet, R. Piarroux, and J. Gaudart (2017). Mathematical models for predicting human mobility in the context of infectious disease spread: Introducing the impedance model. *International Journal of Health Geographics* 16(1), 1–11.
- Saltelli, A. (2019). A short comment on statistical versus mathematical modelling. *Nature communications* 10(1), 3870.
- Saltelli, A., G. Bammer, I. Bruno, E. Charters, M. Di Fiore, E. Didier, W. Nelson Espeland, J. Kay, S. Lo Piano, D. Mayo, R. Pielke, T. Portaluri, T. M. Porter, A. Puy, I. Rafols, J. R. Ravetz, E. Reinert, D. Sarewitz, P. B. Stark, A. Stirling, J. van der Sluijs, and P. Vineis (2020). Five ways to ensure that models serve society: a manifesto. *Nature* 582, 428–484.
- Shumway, R. H. and D. S. Stoffer (2017). *Time Series Analysis and Its Applications: With R Examples*. Springer Texts in Statistics. Cham: Springer International Publishing.
- Simpson, M. J. and O. J. Maclaren (2023). Profile-wise analysis: a profile likelihood-based workflow for identifiability analysis, estimation, and prediction with mechanistic mathematical models. *PLoS Computational Biology* 19(9), e1011515.
- Stocks, T., T. Britton, and M. Höhle (2020). Model selection and parameter estimation for dynamic epidemic models via iterated filtering: Application to rotavirus in Germany. *Biostatistics* 21(3), 400–416.
- Subramanian, R., V. Romeo-Aznar, E. Ionides, C. T. Codeço, and M. Pascual (2020). Predicting re-emergence times of dengue epidemics at low reproductive numbers: DENV1 in Rio de Janeiro, 1986–1990. *Journal of the Royal Society Interface* 17(167), 20200273.
- Toni, T., D. Welch, N. Strelkowa, A. Ipsen, and M. P. Stumpf (2009). Approximate Bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of the Royal Society Interface* 6, 187–202.
- Tracy, M., M. Cerdá, and K. M. Keyes (2018). Agent-based modeling in public health: Current applications and future directions. *Annual Review of Public Health* 39(1), 77–94. PMID: 29328870.
- Trevisin, C., J. C. Lemaitre, L. Mari, D. Pasetto, M. Gatto, and A. Rinaldo (2022). Epidemicity of cholera spread and the fate of infection control measures. *Journal of the Royal Society Interface* 19(188), 20210844.
- Tuite, A. R., J. Tien, M. Eisenberg, D. J. Earn, J. Ma, and D. N. Fisman (2011). Cholera epidemic in Haiti, 2010: Using a transmission model to explain spatial spread of disease and identify optimal control interventions. *Annals of internal medicine* 154(9), 593–601.
- Vandermeer, J. H. and D. E. Goldberg (2013). *Population ecology: first principles*. Princeton University Press.
- Walton, D., A. Suri, and P. Farmer (2011). Cholera in haiti: fully integrating prevention and care. *Annals of internal medicine* 154(9), 635–637.

- Wei, C. H. and F. Zhu (2024). Mean-preserving rounding integer-valued arma models. *Journal of Time Series Analysis* 46, 530–551.
- Wheeler, J. and E. L. Ionides (2023). Likelihood based inference for arma models. *arXiv preprint arXiv:2310.01198*.
- Wheeler, J., N. McAllister, and Sylvertooth (2023). arima2. <https://cran.r-project.org/web/packages/arima2/index.html>.
- Wheeler, J., A. Rosengart, Z. Jiang, K. Tan, N. Treutle, and E. L. Ionides (2024, 04). Informing policy via dynamic models: Cholera in haiti. *PLOS Computational Biology* 20(4), 1–31.
- Whittle, P. (1951). *Hypothesis Testing in Time Series Analysis*, Volume 4. Almqvist & Wiksells boktr.
- Wood, S. N. (2010). Statistical inference for noisy nonlinear ecological dynamic systems. *Nature* 466, 1102–1104.