

Efficient likelihood approximation via gaussian processes:
with an application to an existing *Plasmodium vivax*
malaria model

The University of Melbourne

Jacob Cumming

May 2024

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| I | Literature Review | 3 |
| 2 | Epidemiological Modelling | 5 |
| 2.1 | Deterministic ODE models | 6 |
| 2.2 | Stochastic models | 8 |
| 2.3 | Doob-Gillespie Algorithm | 11 |
| 2.4 | τ -leaping | 11 |
| 3 | Malaria and Malaria Models | 13 |
| 3.1 | Malaria | 13 |
| 3.2 | Malaria Models | 15 |
| 4 | Parameter Inference | 21 |
| 4.1 | Motivation | 21 |
| 4.2 | Frequentist Parameter Estimation | 21 |
| 5 | Parameter Inference | 25 |
| 5.1 | Monte Carlo Integration | 25 |
| 5.2 | Accept-Reject sampling methods | 25 |
| 5.3 | Essentials for MCMC | 26 |
| 5.4 | Motivating the MH/Gibbs Algorithm | 27 |
| 5.5 | Metropolis Hastings | 27 |
| 5.6 | Gibb's Sampling | 29 |
| 5.7 | Diagnostics for Metropolis Hastings | 29 |
| 6 | Gaussian Processes | 31 |
| 6.1 | Motivation and Definitions | 31 |
| 6.2 | Families of Kernel Function | 33 |
| 7 | Gaussian Process Regression | 37 |
| 7.1 | Observation Variance | 40 |
| 7.2 | Model Selection | 41 |
| 7.3 | Differing mean functions | 42 |
| 7.4 | Bayesian Acquisition Functions | 42 |

| | | |
|-----------|---|-----------|
| II | Calibrating Parameters for a <i>P. vivax</i> Model | 45 |
| 8 | Methods | 47 |
| 8.1 | Creation of Synthetic Data | 47 |
| 8.2 | Model Simulations and Discrepancy Function | 48 |
| 8.3 | Gaussian Process and Initialisation | 48 |
| 8.4 | Bayesian Acquisition and Parameter Updates | 49 |
| 9 | Results and Discussion | 51 |
| 9.1 | Results | 51 |
| 9.2 | Discussion | 54 |
| 9.3 | Further Work | 54 |
| | Bibliography | 55 |

List of Tables

| | | |
|-----|--|----|
| 8.1 | Conservative upper bounds for parameters to be calibrated. Values were informed by Champagne et al. 2022; White et al. 2016. All lower bounds were zero. | 49 |
|-----|--|----|

List of Figures

| | | |
|-----|---|----|
| 2.1 | Some simple model schematics, with varying numbers of compartments: S (susceptable), E (exposed), I (infectious) and R (recovered). The force of infection λ_t is usually a function of I_t , depicted by the dashed red lines. μ and ν are natural birth and death rates respectively. γ is the rate of progression out of the infectious state. In each of these models the physical interpretation differs slightly. In the SIS and $SEIR$ models, it is the rate at which individuals move from infectious to susceptible again or into lifelong immunity, whereas in the SI with demography model, it can be interpreted as the increase to the rate of death attributable to disease induced mortality. σ is the rate of progression from a state of latent infection to becoming infectious. | 6 |
| 2.2 | Solutions to the ordinary differential equations describing the models depicted in Figure 2.1. The initial infectious population was $I_0 = 10$, with $S_0 = 990$. In the $SEIR$ model, $E_0 = R_0 = 0$. For all models $\beta = 0.4$. For the SIS and SI model with demography $\gamma = 1/4$. For the SI model with demography $\mu = 0.012$, and $\nu = 0.0012$. For the $SEIR$ model, $\gamma = 1/90$, and $\sigma = 1/2$ | 7 |
| 2.3 | Exact stochastic simulations of the 3 different models using Algorithm 1. The parameters used were identical to those in Figure 2.2 | 12 |
| 3.1 | The <i>P. vivax</i> (malaria) lifecycle. <i>P. falciparum</i> does not have a dormant liver hypnozoite stage. Created with BioRender.com. | 13 |
| 3.2 | A simple Ross-Macdonald malaria model schematic, as described by Aron and May 1982. S_H and I_H are the number of susceptible and infected humans respectively, and S_M and I_M are the number of susceptible and infected mosquitos. The rate of human infection (λ_H) is dependant on I_M , and the rate of human infection (λ_M) is dependant on I_H | 15 |
| 3.3 | Diagram for <i>P. vivax</i> model in a tropical setting described by White et al. 2016. S and I are the number of susceptible and infected humans and mosquitos (denoted by subscript M). $\lambda_H = mabI_M$ and $\lambda_M = ac(I_0 + I_L)$ | 16 |
| 3.4 | Diagram for <i>P. vivax</i> model described by Champagne et al. 2022. $I_{\text{total}} = I_0 + I_L$. Since the mosquito dynamics have been removed, λ now not has no dependencies on the number of infectious mosquitos. | 18 |

| | | |
|-----|---|----|
| 4.1 | Two linear models of the form $f_i(\boldsymbol{\theta}) = \theta_0 + \theta_1 x_i$ fit given the set of observations $\{(1, 2), (2, 4), (3, 4)\}$ using the method of least squares and maximum likelihood under the assumption that the data are independent realisations by a Poisson distribution with $\text{Pois}(f_i(\boldsymbol{\theta}))$. The least squares estimates were $\theta_0^{\text{LSE}} = 4/3$ and $\theta_1^{\text{LSE}} = 1$. The maximum likelihood estimates were $\hat{\theta}_0 \approx 1.329$ and $\hat{\theta}_1 \approx 0.751$ | 22 |
| 5.1 | 30,000 samples from the posterior distribution of p using the Metropolis Hastings algorithm. It was assumed that $p \sim \text{U}(0, 1)$ and $H \sim \text{Binom}(10, p)$, given $H = 6$. A uniform and normal proposal distributions were compared. | 28 |
| 5.2 | Using a basic SIS model, using incident changing likelihood function. | 29 |
| 5.3 | Using Gibbs on beta and gamma given an ' R_0 ' observation | 30 |
| 6.1 | Ten sample realisations from 4 different kernels with hyperparameters $\ell = 1$, and $\sigma_o^2 = 1$. One realisation is bolded. Samples for each kernel were generated from the same seed. | 34 |
| 6.2 | Ten realisations of zero mean Gaussian processes with the squared exponential kernel, varying the length and amplitude parameters. The samples were generated using the same seed | 35 |
| 7.1 | Sequence of Gaussian process regressions on the target function (black) $f(x) = x(x-1)(x+1)$, after 1, 2, 4, and 8 observations in blue. The red lines are new realisations from the conditioned Gaussian process. The Gaussian process was zero mean and had a squared exponential kernel. The hyperparameters were fixed with $\ell = 2.7$ and $\sigma_k^2 = 1.1$ | 39 |
| 7.2 | Sequence of Gaussian process regressions on the target function (black) $f(x) = x(x-1)(x+1)$, after 2, 4, 8, and 16 observations of $f(x_i) + \varepsilon_i$, where ε_i is i.i.d. $\mathcal{N}(0, \sigma_o^2)$ with $\sigma_o^2 = 0.01$ in blue. The red lines are new realisations from the conditioned Gaussian process. The Gaussian process was 0 mean and had a squared exponential kernel. The hyperparameters were fixed with $\ell = 2.7$ and $\sigma_k^2 = 1.1$ | 40 |
| 8.1 | A Doob-Gillespie Simulation of the model described by Champagne et al. 2022 with $\alpha = 0.4$, $\beta = 0.4$, $\gamma_L = 1/223$, $\lambda = 0.04$, $f = 1/72$, $r = 1/60$, and $\delta = 0$. The population was 1000, with 10 initial infections (both blood and liver stage). . . . | 47 |
| 9.1 | Four sample discrepancies $\ln \mathcal{D}(\boldsymbol{\theta})$ taken along 21 values of each parameter to be predicted between the lower and upper parameter bound. The Gaussian process predicts the mean of $\ln \mathcal{D}(\boldsymbol{\theta})$ after initialisation, 50 iterations of Bayesian acquisition, and 500 iterations of Bayesian acquisition. All parameters not in the slice are fixed at the true parameters that generated the synthetic observed data. The black lines are $\mathbb{E}(d_{\mathcal{GP}}(\boldsymbol{\theta}))$ and the red lines are sample realisations from $(d_{\mathcal{GP}}(\boldsymbol{\theta}))$ | 52 |
| 9.2 | Final synthetic likelihoods $\hat{L}(\boldsymbol{\theta})$ after 500 iterations of Bayesian acquisition. Each likelihood shows how that likelihood changes across the parameter space. | 53 |

Chapter 1

Introduction

All the stuff that's really hard about calibrating malaria model parameters. Recent model don't even calibrate using their models.

Why are we looking at mathematical models for vivax malaria and why do we need to do more complicated parameter inference. Outline that the limitation is that it takes a long time to run simulations from complicated malaria models and we still need to fit to data.

The aim of this thesis was to investigate the use of Gaussian Processes to

Part I

Literature Review

Chapter 2

Epidemiological Modelling

In order to study the behaviour and characteristics of disease spread and eradication, compartmental epidemiological models have been developed. They seek to simplify the dynamics of a disease down to a mathematically representable form. Inference on these models allow for an understanding of how the modelled disease spreads, and allows an assessment of how effective differing disease interventions (such as treatments or vaccinations) may be without the need for large long term trials. Models can also simulate various scenarios such as increases or decreases in viral transmission.

Simple compartmental disease models assume individuals can be only be in one of a finite number of states (which are called compartments). These compartments usually correspond to a state of disease. Some simple common compartments include:

- S - Susceptable: at risk of contracting the disease
- E - Exposed: contracted the disease but not yet transmitting it
- I - Infectious (also called Infected): at risk of transmitting the disease
- R - Recovered: neither at risk of contracting or transmitting the disease.

The number of people in each compartment at time t is a (possibly non-deterministic) function of time t , which we indicate as a subscript t (eg. S_t is the number of susceptibles at time t). Models are routinely described by the compartments they contain. For example, an SIS model, is a model with the susceptible and infectious compartments. Recovering from the infection leaves you susceptible to reinfection (for example most sexually transmitted diseases (Keeling and Rohani 2008, p. 56)), and is graphically depicted in Figure 2.1a.

Furthermore we can also include demography into a compartmental model. Diseases that infect the individual until the time of death such as bovine spongiform encephalopathy (BSE commonly known as mad cow disease) may be modelled using an SI model with demography (birth and death rates), as depicted in Figure 2.1b (Hagenaars, Donnelly, and Ferguson 2006).

Childhood diseases such as varicella (chickenpox) which give lifetime immunity after infection can be modelled using an $SEIR$ model (see Figure 2.1c), particularly when modelling a local outbreak setting (for example Zha et al. 2020 used the $SEIR$ model to model a school outbreak of varicella). Not including demography is usually appropriate when disease induced mortality is low.

The number (and names) of compartments can be extended and configured as needed, and compartments could be added for vaccinated individuals, quarantined individuals and so on. By

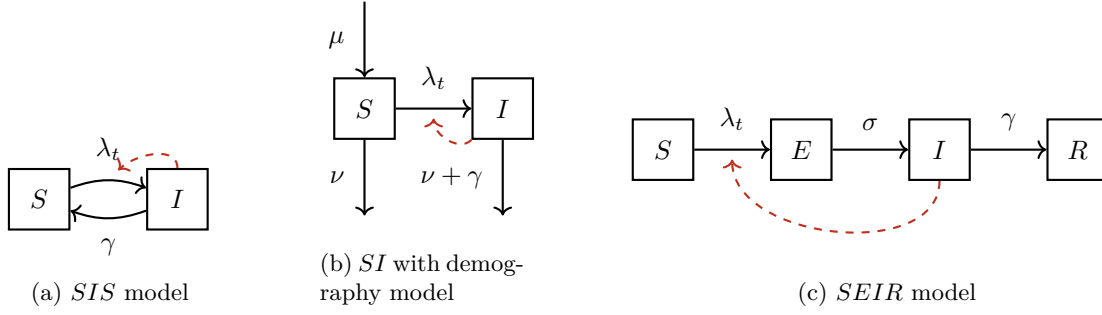


Figure 2.1: Some simple model schematics, with varying numbers of compartments: S (susceptible), E (exposed), I (infectious) and R (recovered). The force of infection λ_t is usually a function of I_t , depicted by the dashed red lines. μ and ν are natural birth and death rates respectively. γ is the rate of progression out of the infectious state. In each of these models the physical interpretation differs slightly. In the *SIS* and *SEIR* models, it is the rate at which individuals move from infectious to susceptible again or into lifelong immunity, whereas in the *SI* with demography model, it can be interpreted as the increase to the rate of death attributable to disease induced mortality. σ is the rate of progression from a state of latent infection to becoming infectious.

convention N_t (often simply N in models with a closed population) is the total number of individuals in the model, the sum of all compartments.

2.1 Deterministic ODE models

Diseases are often simulated as deterministic ordinary differential equations. For the examples below we assume that the force of infection λ_t is proportional to the number of people in I , such that $\lambda_t := \beta \frac{I_t}{N_t}$. β can be interpreted as the average number of people that an individual interacts with per day in a way such that disease would be spread in that interaction per unit of time t . This is sometimes called the effective contact rate. Therefore since $\frac{I_t}{N_t}$ is the probability that a randomly selected individual is infectious, $\beta \frac{I_t}{N_t}$ can be interpreted as the average number of people that a person interacts with each day who are infectious in a way that they would pass on the infection. In different diseases β varies dramatically, as some diseases need prolonged exposure or sexual contact to transmit, whereas some are very highly transmittable, and so will have very low β . Implicitly there is also a (sometimes poor) assumption of complete uniformly random mixing of people. Note that for this thesis we assume that β is frequency dependent (people interact with the same number of people regardless of population size) as opposed to density dependent (people interact with a number of people proportional to population size, in which case $\lambda := \beta I_t$).

The *SIS* model depicted in Figure 2.1a, the ordinary differential equations (ODEs) governing the model could be

$$\frac{dS_t}{dt} = -\lambda S_t + \gamma I_t = -\beta \frac{I_t}{N} S_t + \gamma I \quad (2.1)$$

$$\frac{dI_t}{dt} = \lambda S_t - \gamma I_t = \beta \frac{I_t}{N} S_t - \gamma I. \quad (2.2)$$

With a stated assumption that population size is closed, equation 2.1 fully describes the model.

The system of ODEs that describe the *SI* with demography model is

$$\frac{dS_t}{dt} = \mu N_t - \lambda_t S_t - \nu I_t = \mu N_t - \beta \frac{I_t}{N_t} S_t - \nu I_t \quad (2.3)$$

$$\frac{dI_t}{dt} = \lambda_t S_t - (\gamma + \nu) I_t = \beta \frac{I_t}{N_t} S_t - (\gamma + \nu) I_t. \quad (2.4)$$

Here is it important to note that N_t is not necessarily constant.

Finally, the system of ODEs that describe the *SEIR* model is

$$\frac{dS_t}{dt} = -\lambda_t S_t - \nu I_t = -\beta \frac{I_t}{N} S_t + \gamma I_t \quad (2.5)$$

$$\frac{dE_t}{dt} = \lambda_t S_t - \omega E_t = \beta \frac{I_t}{N} S_t - \omega I_t \quad (2.6)$$

$$\frac{dI_t}{dt} = \omega E_t - \gamma I_t \quad (2.7)$$

$$\frac{dR_t}{dt} = \gamma I_t \quad (2.8)$$

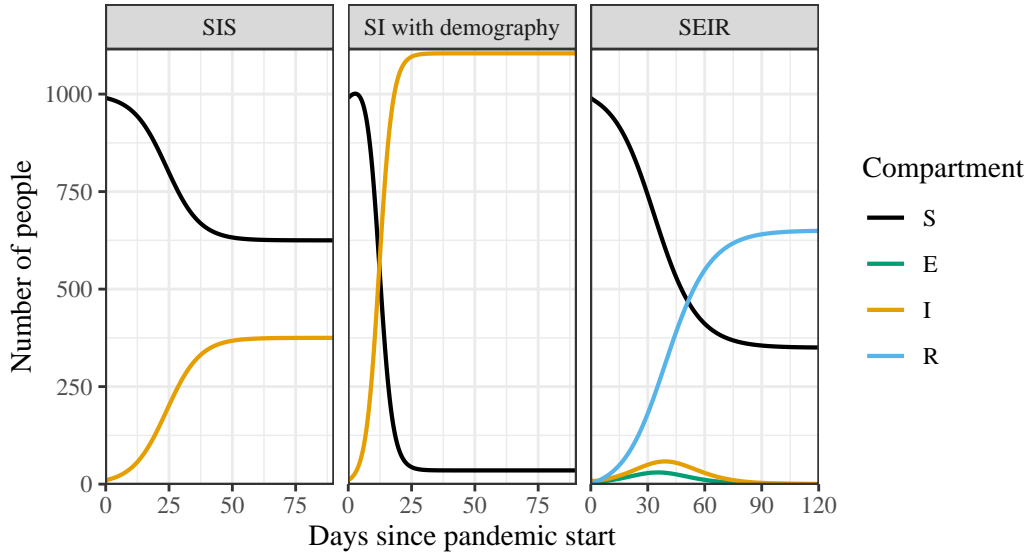


Figure 2.2: Solutions to the ordinary differential equations describing the models depicted in Figure 2.1. The initial infectious population was $I_0 = 10$, with $S_0 = 990$. In the *SEIR* model, $E_0 = R_0 = 0$. For all models $\beta = 0.4$. For the *SIS* and *SI* model with demography $\gamma = 1/4$. For the *SI* model with demography $\mu = 0.012$, and $\nu = 0.0012$. For the *SEIR* model, $\gamma = 1/90$, and $\sigma = 1/2$.

After specifying the initial for each compartment the ordinary differential equations have a deterministic output, such as in 2.2.

2.2 Stochastic models

Motivating the form of the stochastic model

To establish a relationship between the deterministic and stochastic disease models, we first need to establish Poisson point processes and their properties.

Definition 2.1 (Poisson Point Process). $\{\mathcal{N}(t)\}_{t \geq 0}$ is a (stationary) Poisson point process with intensity λ if

1. $\mathcal{N}(0) = 0$
2. $\mathcal{N}(t_1) - \mathcal{N}(t_0), \mathcal{N}(t_2) - \mathcal{N}(t_1), \dots, \mathcal{N}(t_n) - \mathcal{N}(t_{n-1})$ are independent for $0 \leq t_0 < t_1 < \dots < t_{n-1} < t_n$
3. $\mathcal{N}(t_2) - \mathcal{N}(t_1) \sim \text{Pois}(\lambda(t_2 - t_1)), 0 \leq t_1 < t_2$.

Deterministic ODE models are appropriate to study a kind of aggregate disease spread behaviour, and well approximate real world behaviour when the numbers in each compartment are large, however at the start of an epidemic, when the number of infected individuals is small the behaviour of the epidemic may vary significantly. It is possible that if the average number of people that an infectious individual infects near the beginning of the epidemic (formally referred to as R_0) is close to 1, then the disease may die out or become stable. Under the deterministic *SIS* model described by equations 2.1 and 2.2, consider the model at time t^* the instantaneous rate at which S is decreasing is $\beta \frac{I_{t^*}}{N} S_{t^*}$. In other words, one individual leaves the S compartment every $\beta \frac{I_{t^*}}{N} S_{t^*}$ units of time. We can consider a Poisson point process $\{\mathcal{N}_1(t - t^*)\}_{t \geq t^*}$ with intensity $\beta \frac{I_{t^*}}{N} S_{t^*}$ corresponding to the count of the number of individuals who have left S and entered I t units of time since t^* .

$$\frac{d\mathbb{E}(\mathcal{N}_1(t^*))}{dt^*} = \lim_{\delta \rightarrow 0} \frac{\mathbb{E}(\mathcal{N}(t^* + \delta) - \mathcal{N}(t^*))}{\delta} = \frac{\beta \frac{I_{t^*}}{N} S_{t^*} (t^* + \delta - t^*) \beta \frac{I_{t^*}}{N} S_{t^*}}{\delta} = \beta \frac{I_{t^*}}{N} S_{t^*}.$$

Under the same deterministic formulation of the model, the instantaneous rate into S at t^* is γI_{t^*} . Therefore as above we can construct a Poisson point process $\{\mathcal{N}_2(t - t^*)\}_{t \geq t^*}$ with rate γI_{t^*} describing the number of recoveries from I to S , with $\frac{d\mathbb{E}(\mathcal{N}_2(t^*))}{dt^*} = \gamma I_{t^*}$. Combining the two processes, we can see that the rate of change in the average number of people in S is

$$\frac{d\mathbb{E}(\mathcal{N}_2(t^*) - \mathcal{N}_1(t^*))}{dt^*} = \frac{d\mathbb{E}(\mathcal{N}_2(t^*)) - d\mathbb{E}(\mathcal{N}_1(t^*))}{dt^*} = -\beta \frac{I_t}{N} S_t + \gamma I = \frac{dS_t}{dt}.$$

Therefore we can create a stochastic model where the local average in each compartment matches the behaviour of the ODE model at the same state. We do this first by formulating the model as a random vector $\{\mathbf{C}_t\}_{t \geq 0} = \{C_1(t), C_2(t), \dots, C_n(t)\}_{t \geq 0}$ where $C_i : \mathbb{R} \rightarrow \mathbb{N} \cup \{0\}$, is the number of people in compartment C_i , and for any fixed t , $\{C_1(t), C_2(t), \dots, C_n(t)\}$ is a random variable describing the state of the model. For example in a model with S and I compartments, $\{\mathbf{C}_t\}_{t \geq 0} := \{S_t, I_t\}_{t \geq 0}$. $\{\mathbf{C}_t\}_{t \geq 0}$ is a continuous time Markov chain (see Definition 5.3) with transition kernel corresponding to the rates of the model. For example, in the *SI* model with demography in Figure 2.1b the transition rates are:

- $\{s, i\}$ to $\{s + 1, i\}$ has rate $\mu(s + i)$
- $\{s, i\}$ to $\{s - 1, i\}$ has rate νs

- $\{s, i\}$ to $\{s - 1, i + 1\}$ has rate $\beta \frac{i}{i+s} s$
- $\{s, i\}$ to $\{s, i - 1\}$ has rate $(\nu + \gamma)i$.

We can interpret each transition as a separate events, each behaving as independent Poisson point processes until the time of the first transition. Therefore at time t^* we have the Poisson point processes:

- $\{\mathcal{E}_1(t)\}_{t \geq 0}$: the number of births into S after time t^* with intensity μN_{t^*}
- $\{\mathcal{E}_2(t)\}_{t \geq 0}$: the number of deaths in S after time t^* with intensity νS_{t^*}
- $\{\mathcal{E}_3(t)\}_{t \geq 0}$: the number of infections after time t^* with intensity $\beta \frac{I_{t^*}}{N_{t^*}} S_{t^*}$
- $\{\mathcal{E}_4(t)\}_{t \geq 0}$: the number of deaths from I after time t^* with intensity $(\nu + \gamma)I_{t^*}$.

Theorem 2.2 (Sums of Independent Poisson Point Processes). *Given independent Poisson point processes $\{\mathcal{N}_1(t)\}_{t \geq 0}, \{\mathcal{N}_2(t)\}_{t \geq 0}, \dots, \{\mathcal{N}_n(t)\}_{t \geq 0}$, with intensities $\lambda_1, \lambda_2, \dots, \lambda_n$,*

$$\{\mathcal{N}(t)\}_{t \geq 0} := \{\mathcal{N}_1(t) + \mathcal{N}_2(t) + \dots + \mathcal{N}_n(t)\}_{t \geq 0}$$

is a Poisson point process with intensity $\lambda_1 + \lambda_2 + \dots + \lambda_n$.

Proof. We show that $\{\mathcal{N}(t) := \{\mathcal{N}_1(t) + \mathcal{N}_2(t) + \dots + \mathcal{N}_n(t)\}_{t \geq 0}\}$ meets each component of Definition 2.1.

1. $\mathcal{N}(0) := \mathcal{N}_1(0) + \mathcal{N}_2(0) + \dots + \mathcal{N}_n(0) = 0$ since $\mathcal{N}_i(0) = 0$ by definition of a Poisson point process.
2. We show that $\mathcal{N}(t_1) - \mathcal{N}(t_0), \mathcal{N}(t_2) - \mathcal{N}(t_1), \dots, \mathcal{N}(t_n) - \mathcal{N}(t_{n-1})$ with $0 \leq t_0 < t_1 < \dots < t_n$ are independent.

$$\mathcal{N}(t_i) - \mathcal{N}(t_{i-1}) = \underbrace{[\mathcal{N}_1(t_i) - \mathcal{N}_1(t_{i-1})]}_{X_{i1}} + \underbrace{[\mathcal{N}_2(t_i) - \mathcal{N}_2(t_{i-1})]}_{X_{i2}} + \dots + \underbrace{[\mathcal{N}_n(t_i) - \mathcal{N}_n(t_{i-1})]}_{X_{in}}$$

X_{ik} is independent of $X_{j\ell}$ for $k \neq \ell$ since \mathcal{N}_k and \mathcal{N}_ℓ are independent processes. X_{ik} is independent of X_{jk} for $i \neq j$ by the second property of Definition 2.1. Therefore all X_{ik} are independent of $X_{j\ell}$ for $i \neq j$, and all j, k . Hence $\mathcal{N}(t_1) - \mathcal{N}(t_0), \mathcal{N}(t_2) - \mathcal{N}(t_1), \dots, \mathcal{N}(t_n) - \mathcal{N}(t_{n-1})$ with $0 \leq t_0 < t_1 < \dots < t_n$ are independent.

3. For fixed $t_1 < t_2$, and $i \in \{1, 2, \dots, n\}$,

$$\mathcal{N}_i(t_2) - \mathcal{N}_i(t_1) \sim \text{Pois}((t_2 - t_1)\lambda_i).$$

Consider the associated moment generating function of $\mathcal{N}_i(t_2) - \mathcal{N}_i(t_1)$,

$$M_i(z) := \exp(\lambda_i(t_2 - t_1)(\exp(z) - 1)).$$

Therefore the moment generating function of

$$\mathcal{N}(t_2) - \mathcal{N}(t_1) = [\mathcal{N}_1(t_2) - \mathcal{N}_1(t_1)] + [\mathcal{N}_2(t_2) - \mathcal{N}_2(t_1)] + \dots + [\mathcal{N}_n(t_2) - \mathcal{N}_n(t_1)]$$

is

$$M(z) := \prod_{i=1}^n M_i(z) = \exp[(\lambda_1(t_2 - t_1) + \lambda_2(t_2 - t_1) + \cdots + \lambda_n(t_2 - t_1))(\exp(z) - 1)].$$

Therefore $\mathcal{N}_1(t) + \mathcal{N}_2(t) + \cdots + \mathcal{N}_n(t) \sim \text{Pois}((\lambda_1 + \lambda_2 + \cdots + \lambda_n)t)$ by the uniqueness of the moment generating function.

□

Theorem 2.3 (Time to First Event in Poisson Point Process). *Given a Poisson point process $\{\mathcal{N}(t)\}_{t \geq 0}$ with intensity λ , let $\tau = \inf\{t | \mathcal{N}(t_0 + t) - \mathcal{N}(t_0) = 1, t > 0\}$. $\tau \sim \text{Exp}(\lambda)$ for $t_0 \geq 0$*

Proof.

$$\Pr(\tau > x) = \Pr(\mathcal{N}(t_0 + x) - \mathcal{N}(t_0) = 0) = \frac{(\lambda x)^0 e^{-\lambda x}}{0!} = e^{-\lambda x}$$

□

By Theorem 2.2 and Theorem 2.3,

$$\{\mathcal{E}(t)\}_{t \geq 0} := \{\mathcal{E}_1(t) + \mathcal{E}_2(t) + \mathcal{E}_3(t) + \mathcal{E}_4(t)\}_{t \geq 0}$$

is a Poisson point process with intensity

$$\mu N_{t^*} + \nu S_{t^*} + \beta \frac{I_{t^*}}{N_{t^*}} S_{t^*} + (\nu + \gamma) I_{t^*},$$

and the time to the next event is random variable distributed

$$\text{Exp}(\mu N_{t^*} + \nu S_{t^*} + \beta \frac{I_{t^*}}{N_{t^*}} S_{t^*} + (\nu + \gamma) I_{t^*}).$$

Theorem 2.4 (Probability of i th Poisson Process Generating the Next Event). *Consider independent Poisson point processes*

$$\{\mathcal{N}_1(t)\}_{t \geq 0}, \{\mathcal{N}_2(t)\}_{t \geq 0}, \dots, \{\mathcal{N}_n(t)\}_{t \geq 0}$$

having intensities $\lambda_1, \lambda_2, \dots, \lambda_n$. For fixed t_0 , let $\tau_i := \inf\{t | \mathcal{N}(t_0 + t) - \mathcal{N}(t_0) = 1\}$. Then

$$\Pr(\min_i \tau_i = \tau_j) = \frac{\lambda_j}{\sum_{i=1}^n \lambda_i}.$$

Proof. By Theorem 2.3, $\tau_i \sim \text{Exp}(\lambda_i)$. Therefore

$$\begin{aligned}
\Pr(\min_i \tau_i = \tau_j) &= \int_0^\infty \Pr(\{\tau_i = x\} \cup \bigcup_{j \neq i} \{\tau_j > x\}) dx \\
&= \int_0^\infty \Pr(\{\tau_i = x\} \cup \bigcup_{j \neq i} \{\tau_j > x\}) dx \\
&= \int_0^\infty \Pr(\tau_i = x) \times \prod_{j \neq i} \Pr(\tau_j > x) dx && \text{(by independence)} \\
&= \int_0^\infty \lambda_i \exp(-\lambda_i x) \times \prod_{j \neq i} \exp(-\lambda_j x) dx \\
&= \lambda_i \int_0^\infty \exp(-(\sum_{i=1}^n \lambda_j)x) dx \\
&= \lambda_i \left[\frac{\exp(-(\sum_{i=1}^n \lambda_j)x)}{\sum_{i=1}^n \lambda_j} \right]_0^\infty \\
&= \frac{\lambda_i}{\sum_{i=1}^n \lambda_j}
\end{aligned}$$

□

2.3 Doob-Gillespie Algorithm

All of this leads naturally to a common method of simulating the stochastic model. The Doob-Gillespie algorithm (often just called the Gillespie algorithm) is an algorithm that simulates a stochastic realisation of a model given a set of starting conditions.

Algorithm 1 The Doob-Gillespie Algorithm

- 1: Initialise time $t \leftarrow 0$ and initial state of the model $\mathbf{C}(0) := \{C_1(0), C_2(0), \dots, C_n(0)\}$
 - 2: **while** termination condition not met **do**
 - 3: Calculate intensities λ_i for all possible events \mathcal{E}_i
 - 4: Calculate total intensity $\lambda = \sum_i \lambda_i$
 - 5: Generate $\Delta t \sim \text{Exp}(\lambda)$
 - 6: Choose event E_i with probability $\frac{\lambda_i}{\lambda}$
 - 7: Update time $t \leftarrow t + \Delta t$
 - 8: Update state of $\mathbf{C}(t + \delta t) \leftarrow \mathbf{C}(t) +$ change in state due to event \mathcal{E}_i
 - 9: **end while**
-

2.4 τ -leaping

τ -leaping exploits the local Poisson point process like behaviour of epidemiological models. Consider the SIS model, when $S_t = I_t = 10000$. Events happen at a very high rate, meaning the Δt found in each step of the Doob-Gillespie algorithm will be very small, but the rates also change a negligible amount after each event (compare $\gamma \times 10000$ to $\gamma \times 10001$ or $\gamma \times 9999$). Therefore we can approximate the number of events in a short time period τ as a Poisson point process with the total intensity $\lambda = \sum_i \lambda_i$ at time t , with the probability of any one event having the same probability as above of $\frac{\lambda_i}{\lambda}$. Therefore we have the following algorithm.

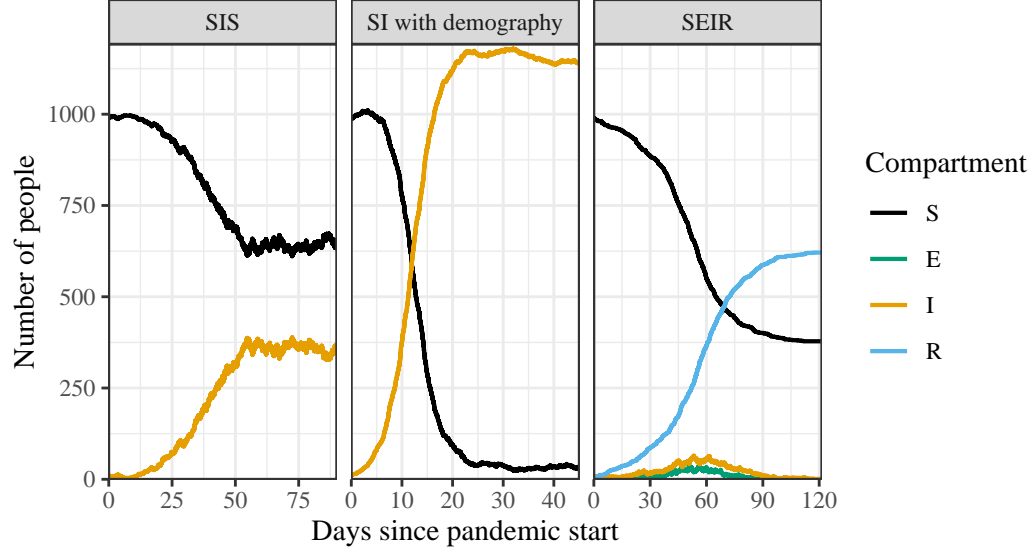


Figure 2.3: Exact stochastic simulations of the 3 different models using Algorithm 1. The parameters used were identical to those in Figure 2.2

Algorithm 2 τ -Leaping Algorithm

- 1: Initialise time $t \leftarrow 0$ and initial state of the model $\mathbf{C}(0) := \{C_1(0), C_2(0), \dots, C_n(0)\}$
 - 2: **while** termination condition not met **do**
 - 3: Calculate intensities λ_i for all possible events \mathcal{E}_i
 - 4: Calculate total intensity $\lambda = \sum_i \lambda_i$
 - 5: Choose a suitable time step τ (this can be deterministic or adaptive)
 - 6: Calculate Poisson random variable $X \sim \text{Poisson}(\lambda)$
 - 7: **for** i in 1 to X **do**
 - 8: Choose event E_i with probability $\frac{\lambda_i}{\lambda}$
 - 9: Update state of $\mathbf{C}(t + \tau) \leftarrow \mathbf{C}(t) + \text{change in state due to event } \mathcal{E}_i$
 - 10: **end for**
 - 11: Update time $t \leftarrow t + \tau$
 - 12: **end while**
-

Chapter 3

Malaria and Malaria Models

3.1 Malaria

Malaria kills around 600,000 people each year, with over 75% of deaths occurring in children under 5 years old (World Health Organization 2022).

Plasmodium Life Cycle

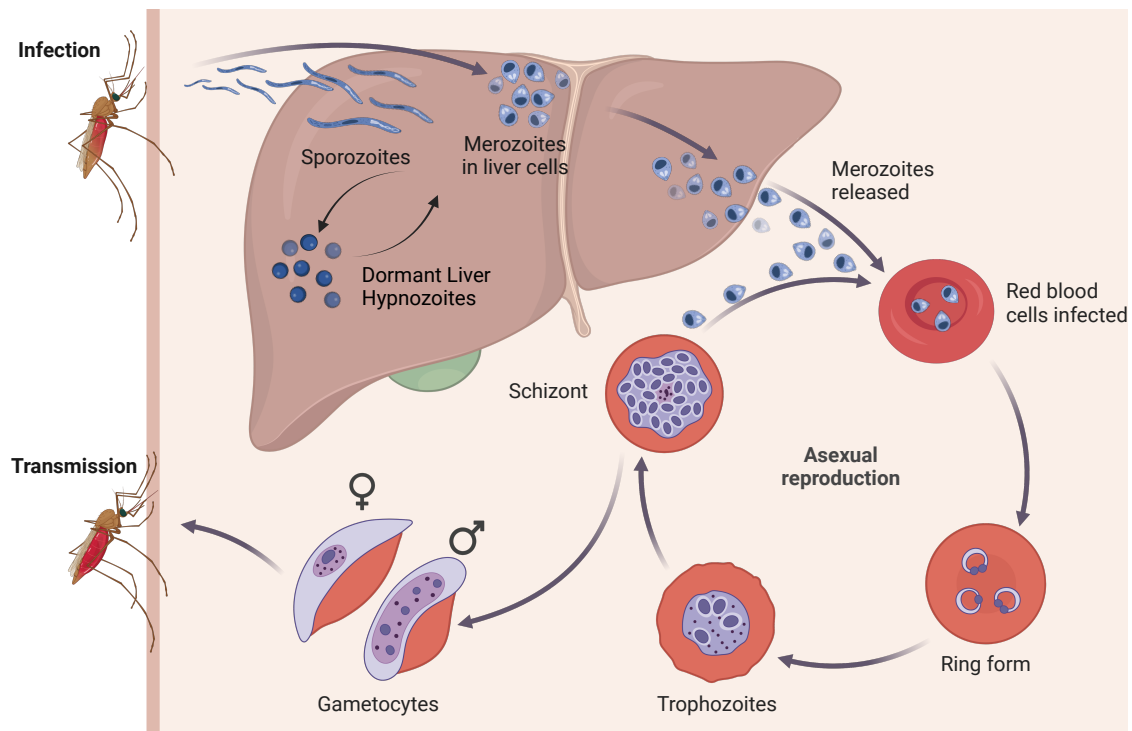


Figure 3.1: The *P. vivax* (malaria) lifecycle. *P. falciparum* does not have a dormant liver hypnozoite stage. Created with BioRender.com.

Malaria is a vector borne disease, needing both human (or other vertebrate) and mosquito hosts to complete its lifecycle (Figure 3.1). Six species of the unicellular parasite are able to infect humans (Milner 2018). Although *Plasmodium falciparum* is responsible for around 90% of total human

malaria deaths, outside of Africa *Plasmodium vivax* is the leading cause of malaria infection (Zekar and Sharman 2023; Adams and Mueller 2017). Sporozoites (a stage of the malaria parasite) enter the human blood stream via the skin after the female mosquito has a blood meal. From the blood stream they proceed to enter into the liver. Once a hepatocyte (liver cell) is invaded, the parasite will undergo asexual reproduction into up to 40,000 merozoites per hepatocyte, which are released into the blood stream. These merozoites then bind to, and invade erythrocytes (red blood cells), once again reproducing 16-32 fold in a process called schizogony. At this point, the erythrocyte membrane is ruptured, allowing for *Plasmodium* to invade new erythrocytes. Eventually, the merozoites undergo sexual differentiation, resulting in the sequestration and maturation of male and female gametocytes in the bone marrow, until they are released into the blood stream to be consumed by a mosquito during a blood feed where it matures into sporozoites ready to reinfect a new vertebrate host when the mosquito next takes a blood feed (Cowman et al. 2016).

Illness, Treatment, and Immunity

The most common symptom of malaria infection in persons without natural or acquired immunity is fever. After treatment, fever will usually subside over a few days. In severe cases, malaria can lead to anemia, cerebral malaria (coma), and respiratory distress (Cowman et al. 2016).

In a population with stable malarial infection, immunity increases with age, with the proportion of severe cases negligible after age 10, and asymptomatic infection being the dominant infection type beyond age 15 (Cowman et al. 2016).

Control and Eradication

Widespread use of DDTs in the mid 20th century led to significant successes in some countries towards the control and eradication of malaria. In the 1980s and 1990s, drug resistant malaria led to a doubling of malaria-attributable death. Currently, the control techniques include insecticide treated bed nets, and a mixture of antimalarial drugs (Cowman et al. 2016).

Plasmodium vivax

Unlike *P. falciparum*, *P. vivax* has hypnozoites, which are a dormant liver stage of the parasite. These can remain dormant for weeks and even months, leading to recurrent infections and illness, possibly until the conditions for transmission are more favourable. In subtropical/temperate areas, the incubation periods can be between 8-12 months, compared to 3-4 weeks in tropical regions. Price et al. 2020. *P. vivax* also has lower levels of the blood stage parasite during infection, which means diagnosis is more difficult, and it has an increased proportion of asymptomatic cases (Adams and Mueller 2017).

It is likely that death and severe disease attributable to *P. vivax* has been traditionally underestimated. In view of recent evidence, the old notion that *P. vivax* is benign has become untenable (Cowman et al. 2016).

Motivating Malaria Models

Levels of asymptomatic cases and latent parasite (in the case of *P. vivax*) are impossible or difficult to experimentally determine without mass testing. By creating a model of the disease, and calibrating the model so that it simulates symptomatic case levels reported by health authorities,

it is possible to estimate these previously ‘hidden’ levels. Furthermore, modelling malaria allows for modelling the effect of public health interventions such as mass treatment or testing, in order to determine an estimate of how effective the intervention may be, before large amounts of money are spent on trials.

3.2 Malaria Models

Ross-Macdonald

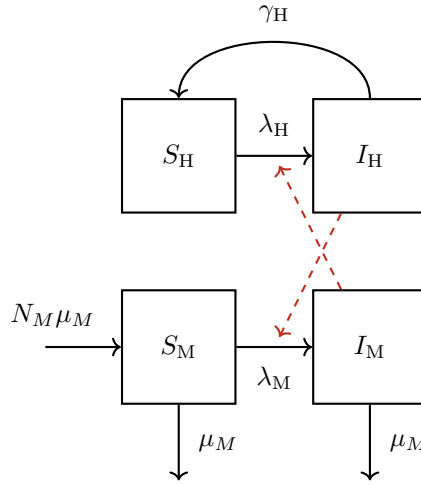


Figure 3.2: A simple Ross-Macdonald malaria model schematic, as described by Aron and May 1982. S_H and I_H are the number of susceptible and infected humans respectively, and S_M and I_M are the number of susceptible and infected mosquitos. The rate of human infection (λ_H) is dependant on I_M , and the rate of human infection (λ_M) is dependant on I_H .

Modelling malaria presents an additional challenge, as the disease is transmitted from mosquito to human and human to mosquito, rather than having direct human to human transmission. The most simple Ross-Macdonald model is depicted in figure 3.2. The ODEs for this model are

$$\begin{aligned}\frac{dS_H}{dt} &= \gamma_H I_H - b T_{HM} I_M \frac{S_H}{N_H} \\ \frac{dI_H}{dt} &= b T_{HM} I_M \frac{S_H}{N_H} - \gamma_H I_H \\ \frac{dS_M}{dt} &= N_M \mu_M + \gamma_M I_M - b T_{MH} S_M \frac{I_H}{N_H} - S_M \mu_M \\ \frac{dI_M}{dt} &= b T_{MH} S_M \frac{I_H}{N_H} - \gamma_M I_M\end{aligned}$$

where b is the biting rate per mosquito, and T_{HM} is the probability of tranmission to a human given a bite by an infectious mosquito, with T_{MH} being vice-versa. Note that it is $\frac{I_H}{N_H}$ in the mosquito dynamics. Biologically this is assuming the number of blood meals a mosquito takes per day is invariant to the size of the human population. Mosquitos don’t ‘recover’ from malaria due to their short lifespans, but the births and deaths are mathematically equivalent to assuming that the rate of ‘recovery’ amongst mosquitos is $\mu_M I_M$ per unit time, with no population dynamics.

A Ross-Macdonald style model simplifies the lifecycle of malaria to the following four steps (Smith et al. 2012):

1. Malaria is transmitted to human (or vertebrate) via a blood feed.
2. Malaria proliferates in the human host until it circulates in the peripheral blood
3. A mosquito then takes a blood feed, ingesting the pathogen
4. Malaria develops within the mosquito host, progressing to its salivary glands, able to infect a human.

Models of *P. Vivax* Malaria

White Model

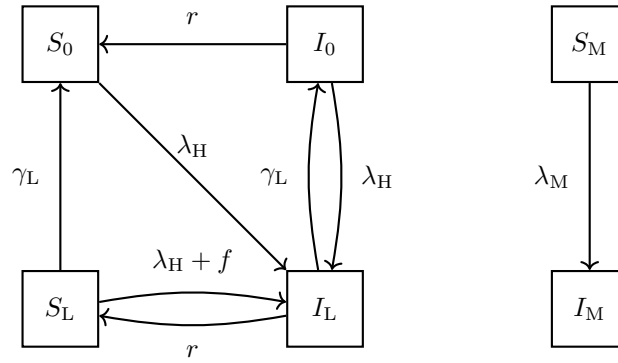


Figure 3.3: Diagram for *P. vivax* model in a tropical setting described by White et al. 2016. S and I are the number of susceptible and infected humans and mosquitos (denoted by subscript M). $\lambda_H = mabI_M$ and $\lambda_M = ac(I_0 + I_L)$

The White model - described in White et al. 2016 (tropical model) and depicted in Figure 3.3 - is characterised by the following ordinary differential equations:

$$\frac{dS_0}{dt} = -\lambda_H S_0 + r I_0 + \gamma_L S_L$$

$$\frac{dI_0}{dt} = -\lambda_H I_0 - r I_0 + \gamma_L I_L$$

$$\frac{dS_L}{dt} = -\lambda_H S_L + r I_L - f S_L - \gamma_L S_L$$

$$\frac{dI_L}{dt} = \lambda_H (S_0 + I_0 + S_L) - r I_L + f S_L - \gamma_L I_L$$

$$\frac{dS_M}{dt} = g - \lambda_M (p S_M - (1-p) I_M) - g S_M$$

$$\frac{dI_M}{dt} = \lambda_M (p S_M - (1-p) I_M) - g I_M. \quad (I_0 + I_L = \text{total number of bloodstage infections})$$

It modifies the Ross-Macdonald models, to capture the differences in disease progression between *P. vivax* and *P. falciparum*. In particular, the White model includes the dormant liver stage that is unique to *P. vivax*.

The model is comprised of six compartments:

1. **S_0 (Susceptible Individuals - No Latent Hypnozoite Liver Stage Infection) :** People in this compartment have no form of malarial infection. These people are susceptible to new malarial infections, and are infected into compartment I_L (with both blood and liver stage parasites) with rate λ_H .
2. **I_L (Infected Individuals - Both Blood Stage and Latent Hypnozoite Liver Stage Infection) :** Individuals in this compartment have both an active blood-stage infection, and latent hypnozoite infection in the liver. They can progress to either I_0 through the clearance of liver stage infection with rate γ_L , or to S_L through clearance of blood stage infection with rate r .
3. **I_0 (Infected Individuals - Blood-Stage Infection Only) :** Those in this compartment have a blood-stage infection with no latent hypnozoite infection in the liver. They are be reinfected into I_L with rate λ_H , relapse with rate f . Blood-stage infection is cleared (moving into compartment S_0) with rate r .
4. **S_L (Susceptable Individuals - Blood-Stage Infection Only) :** Those in this compartment have latent hypnozoite infection in the liver without blood-stage infection. They get novel infection through a mosquito bite into I_L with rate λ_H , or hypnozoite activation with rate f . This means that those in S_L move to compartment I_L with total rate $\lambda_H + f$. Alternatively the hypnozoites are cleared from the liver (moving to compartment S_0) with rate γ_L .
5. **S_M (Susceptable Mosquitoes) :** Suspectable mosquitoes become infectious at rate $\lambda_M p$. They die at rate $g + \lambda_M(1 - p)$. Since there is a constant mosquito population assumption, mosquitoes are born into this state at rate $g + \lambda_M$.
6. **I_M (Infectious Mosquitoes):** Infectious mosquitos die at rate $g + \lambda_M(1 - p)$.

$\lambda_H := mabI_M$ where m is the number of mosquitos per human (held constant since there is no birth or death in the human dynamics), a is the mosquito biting rate, and b is the probability that a human bitten by an infectious mosquito develops an infection.

$\lambda_M := ac(I_0 + I_L)$ where a is defined above, and c is the probability that a mosquito bite on an infectious mosquito causes the mosquito to become infectious. g can be interpreted as the natural birth/death rate for mosquitos. p is then the proportion of mosquitos that survive long enough after the initial infection that the parasite matures enough in the mosquito before becoming infectious to new susceptible humans. Under the assumption that time until parasite transmissability after infection in a mosquito is a constant n days, and that mosquitoes naturally die at rate g , $p = e^{-gn}$. To see this let $V \sim \text{Exp}(g)$, represent the lifespan of the mosquito. $\Pr(V > n) = 1 - F_V(n) = 1 - (1 - e^{-gn}) = e^{-gn}$.

$\lambda_M(1 - p)$ can be interpreted as an additional rate of death, where of the mosquitos that would develop malaria after a bite, a proportion $1 - p$ die instantly. This applies to both the susceptible and infectious mosquitoes. Presumably this approximates a model where mosquitoes are moved to an ‘exposed’ compartment for n time, after initial infection, however no justification is given in (White et al. 2016) for this additional parameter n . A more straightforward SI model could be constructed that absorbs c and n into the single parameter c^* , such that it becomes the proportion of mosquito bites on blood stage infectious humans that result in mosquito infection where the

mosquito does not die before becoming infectious. With steady mosquito population, the mosquito dynamics would now be characterised by

$$\frac{dI_M}{dt} = \lambda_M^* S_M - g I_M \quad \text{where } \lambda_M^* := ac^*(I_0 + I_L).$$

By modelling both liver and bloodstage infection, blood stage infections from relapses can be captured in the dynamics, meaning it is possible to analyse case number data that may be confounded by relapses as well as novel infections.

This model does not account for continual depletion of liver stage parasites which would vary the rate of relapse over time (through clearance or relapse). It also does not directly model any interventions or case importations. The lack of population dynamics means the model may only be useful on a small time scale. Finally, it doesn't account for any importation of disease from an outside area, so if $S_0 = 1$, *P. vivax* is presumed permanently eradicated.

Champagne Model

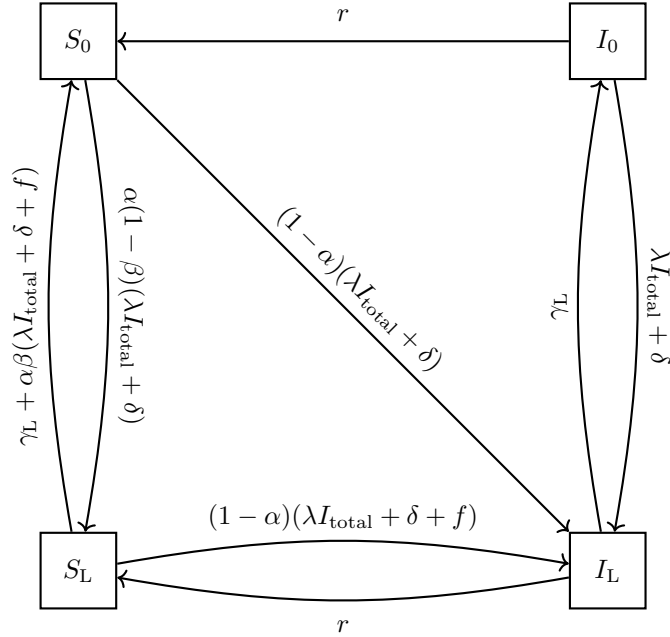


Figure 3.4: Diagram for *P. vivax* model described by Champagne et al. 2022. $I_{\text{total}} = I_0 + I_L$. Since the mosquito dynamics have been removed, λ now not has no dependencies on the number of infectious mosquitoes.

The Champagne model - described in (Champagne et al. 2022) and diagrammatically depicted in figure 3.4 - both simplifies and extends the White model. The model assumes human to human transmission, removing mosquito dynamics, and extends it by adding in a rate of imported cases and treatment of malarial infection. It is characterised by the system of ordinary differential

equations

$$\begin{aligned}
\frac{dI_L}{dt} &= (1 - \alpha)(\lambda I_{\text{total}} + \delta)(S_0 + S_L) + (\lambda I_{\text{total}} + \delta)I_0 + (1 - \alpha)fS_L - \gamma_L I_L - rI_L \\
\frac{dI_0}{dt} &= -(\lambda I_{\text{total}} + \delta)I_0 + \gamma_L I_L - rI_0 \\
\frac{dS_L}{dt} &= -(1 - \alpha(1 - \beta))(\lambda I_{\text{total}} + \delta + f)S_L + \alpha(1 - \beta)(\lambda I_{\text{total}} + \delta)S_0 - \gamma_L S_L + rI_L \\
\frac{dS_0}{dt} &= -(1 - \alpha\beta)(\lambda I_{\text{total}} + \delta)S_0 + (\lambda I_{\text{total}} + \delta)\alpha\beta S_L + \alpha\beta fS_L + \gamma_L S_L + rI_0
\end{aligned}$$

where $I_{\text{total}} := I_0 + I_L$.

The compartments have the same interpretation as in the White model, however the rates between compartments are significantly modified.

The new parameters are λ : the rate of infection, δ : importation rate, α : proportion of those infected who clear blood stage infections through immediate treatment, and β : proportion of those cleared of blood stage infection who are also cleared of liver stage parasites (radical cure) In other words, the proportion of infected individuals $\alpha\beta$ are completely cured from liver and blood stage parasites. The model assumes treatment clears infection instantaneously. Individuals in S_L who relapse or get a new infection are assumed to be cured with the same proportions as new infections from S_0 , but individuals in I_0 who are superinfected are assumed not to seek treatment.

In contrast to the White model, the Champagne model allows analysis of potential treatment interventions, or how much of an impact limiting the importation rate might have on case numbers (through border control/testing). Although the lack of mosquito dynamics simplifies the model and it's running, it is unrealistic. The model still has some of the same problems as the White model, such as not incorporating hypnozoite depletion rates and a lack of population dynamics, meaning all analytic results are done assuming the system is at equilibrium.

Chapter 4

Parameter Inference

4.1 Motivation

Building mathematical models of real world phenomenon allows for us to simulate changes in the world without having to undertake large scale experiments. However, once we have a model that sufficiently approximates *P. vivax* transmission or anything else we are trying to model, we then need to estimate what the ‘true’ underlying parameters are. To do this we calibrate the model against real world data such as case counts, and prevalence surveys. Under frequentist assumptions, there is a ‘true’ set of parameters that if used in our model, simulated the observed data. Under a Bayesian assumption, the parameters are considered to be random, and This chapter explores statistical inference techniques to recover the parameters, under both the frequentist and Bayesian frameworks.

4.2 Frequentist Parameter Estimation

Assume the model is parametrised by a set of parameters $\theta \in \Theta$ which we are trying to estimate by considering some observed data $\mathbf{y}^{\text{obs}} = (y_1^{\text{obs}}, \dots, y_n^{\text{obs}})$. Consider $\mathbf{f}(\theta) = (f_1(\theta), \dots, f_n(\theta))$ some model simulation of \mathbf{y}^{obs} . Often the observed data has some underlying index set x_1, \dots, x_n , where x_i might be something like time. In this case we can also consider the observed data to be $\{(x_1, y_1^{\text{obs}}), \dots, (x_n, y_n^{\text{obs}})\}$, and the model simulated data to be $\{(x_1, f_1(\theta)), \dots, (x_n, f_n(\theta))\}$.

Least Squares Estimator

It is common that models are not random, but instead model the mean behaviour of a system. In this case, $\mathbf{f}(\theta)$ is not random. Therefore we can assume that $y_i^{\text{obs}} = f_i(\theta) + \varepsilon_i$, where ε_i is a random variable with some (possibly unknown) distribution, and zero mean.

When the distribution of ε_i is unknown, a common approach for estimating θ^{LSE} is to take the least squares estimate.

Definition 4.1 (Least Squares Estimate). *The least squares estimate θ^{LSE} for θ is*

$$\theta^{\text{LSE}} := \arg \min_{\theta \in \Theta} \sum_{i=1}^n (f_i(\theta) - y_i^{\text{obs}})^2.$$

Example 4.2. Consider the observed data $\{(x_1, y_1^{\text{obs}}), (x_2, y_2^{\text{obs}}), (x_3, y_3^{\text{obs}})\} = \{(1, 2), (2, 4), (3, 4)\}$, which we assume were generated from the model $f_i(\boldsymbol{\theta}) + \varepsilon_i$, where $f_i(\boldsymbol{\theta}) = \theta_0 + \theta_1 x_i$, and $\mathbb{E}(\varepsilon_i) = 0$. We derive the least squares estimate of our parameters $\boldsymbol{\theta} = (\theta_0, \theta_1)$ by

$$\begin{aligned}\boldsymbol{\theta}^{\text{LSE}} &= \arg \min_{\boldsymbol{\theta}} \left[\sum_{i=1}^3 (f_i(\boldsymbol{\theta}) - y_i^{\text{obs}})^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} \left[\sum_{i=1}^3 (\theta_0 + \theta_1 x_i - y_i^{\text{obs}})^2 \right] \\ &= \arg \min_{\boldsymbol{\theta}} [(\theta_0 + \theta_1 - 2)^2 + (\theta_0 + 2\theta_1 - 4)^2 + (\theta_0 + 3\theta_1 - 4)^2]\end{aligned}$$

Since the expanded quadratic will have positive coefficients out the front of θ_0 and θ_1 , we can solve for $\boldsymbol{\theta}^{\text{LSE}}$ by

$$\begin{aligned}0 &= \frac{\partial}{\partial \boldsymbol{\theta}} [(\theta_0^{\text{LSE}} + \theta_1^{\text{LSE}} - 2)^2 + (\theta_0^{\text{LSE}} + 2\theta_1^{\text{LSE}} - 4)^2 + (\theta_0^{\text{LSE}} + 3\theta_1^{\text{LSE}} - 4)^2] \\ &= \begin{bmatrix} 6\theta_0^{\text{LSE}} + 12\theta_1^{\text{LSE}} - 20 \\ 12\theta_0^{\text{LSE}} + 28\theta_1^{\text{LSE}} - 44 \end{bmatrix}\end{aligned}$$

And solving these two equations results in $\theta_0^{\text{LSE}} = 4/3$ and $\theta_1^{\text{LSE}} = 1$. This can be visually seen in Figure 4.1

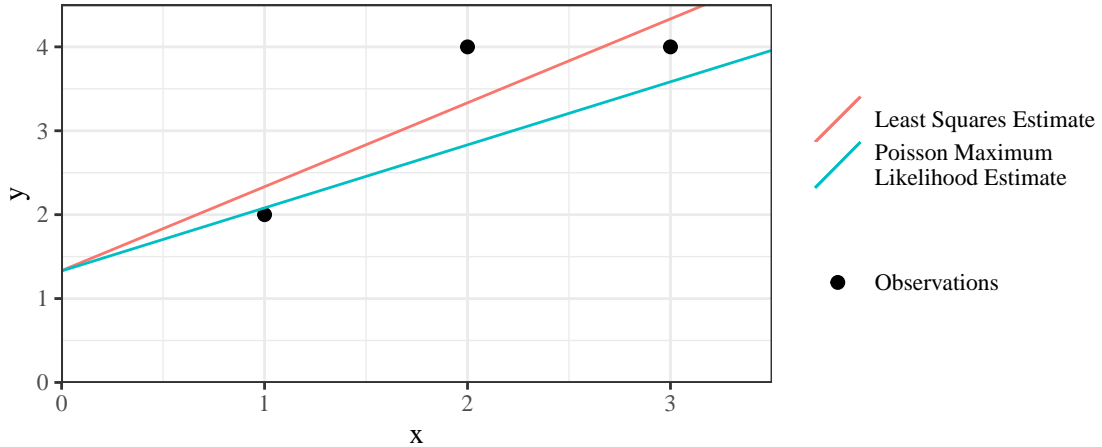


Figure 4.1: Two linear models of the form $f_i(\boldsymbol{\theta}) = \theta_0 + \theta_1 x_i$ fit given the set of observations $\{(1, 2), (2, 4), (3, 4)\}$ using the method of least squares and maximum likelihood under the assumption that the data are independent realisations by a Poisson distribution with $\text{Pois}(f_i(\boldsymbol{\theta}))$. The least squares estimates were $\theta_0^{\text{LSE}} = 4/3$ and $\theta_1^{\text{LSE}} = 1$. The maximum likelihood estimates were $\hat{\theta}_0 \approx 1.329$ and $\hat{\theta}_1 \approx 0.751$.

Maximum Likelihood Estimator

The least square method makes no explicit assumptions about the distribution of the noise ε . However if the distribution of ε is known (or can be reasonably assumed), we can explicitly calculate the probability of the data given the parameters.

Definition 4.3 (Likelihood function). *With \mathbf{y}^{obs} fixed, the likelihood function is*

$$\mathcal{L}(\boldsymbol{\theta}) := \Pr(\mathbf{f}(\boldsymbol{\theta}) + \varepsilon = \mathbf{y}^{\text{obs}} | \boldsymbol{\theta}).$$

Particularly, if $f_i(\boldsymbol{\theta}) + \varepsilon_i$ are independent

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{i=1}^n \Pr(f_i(\boldsymbol{\theta}) + \varepsilon_i = y_i^{\text{obs}} | \boldsymbol{\theta}).$$

In the continuous case, if $\mathbf{f}(\boldsymbol{\theta}) + \varepsilon$ has joint density g ,

$$\mathcal{L}(\boldsymbol{\theta}) := g(\mathbf{y}^{\text{obs}} | \boldsymbol{\theta}).$$

Note that the dependence on \mathbf{y}^{obs} is suppressed, but can be explicitly written as $\mathcal{L}(\boldsymbol{\theta} | \mathbf{y}^{\text{obs}})$. A natural estimate for $\boldsymbol{\theta}$ is the one that maximises the likelihood function, as it coincides with the value of $\boldsymbol{\theta}$ maximises the probability of the data. Such an estimate is called the maximum likelihood estimate.

Definition 4.4 (Maximum Likelihood Estimate). *The maximum likelihood estimate of $\boldsymbol{\theta}$ is*

$$\hat{\boldsymbol{\theta}} := \arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta})$$

It is often computationally easier to deal with the log-likelihood $\ell(\boldsymbol{\theta}) := \ln \mathcal{L}(\boldsymbol{\theta})$. Since \ln is a monotonic function, $\arg \max_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Theta} \ell(\boldsymbol{\theta})$

Example 4.5. *Using the same observed data set as Example 4.2, we assume that y_i^{obs} were generated independently from $f_i(\boldsymbol{\theta}) + \varepsilon_i \sim \text{Pois}(f_i(\boldsymbol{\theta}))$, where $f_i(\boldsymbol{\theta}) = \theta_0 + \theta_1 x_i$ as previously defined. Therefore the maximum likelihood estimate of $\boldsymbol{\theta}$ is*

$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \arg \max \ell(\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^3 y_i^{\text{obs}} \ln(f_i(\boldsymbol{\theta})) - y_i^{\text{obs}}(\boldsymbol{\theta}) - \ln(y_i^{\text{obs}}!) \\ &= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^3 y_i^{\text{obs}} \ln(\theta_0 + \theta_1 x_i) - \theta_0 - \theta_1 x_i - \ln(y_i^{\text{obs}}!) \\ &= \arg \max_{\boldsymbol{\theta} \in \Theta} 2 \ln(\theta_0 + \theta_1) - \theta_0 - \theta_1 + 4 \ln(\theta_0 + 2\theta_1) - \theta_0 - 2\theta_1 + 4 \ln(\theta_0 + 3\theta_1) - \theta_0 - 3\theta_1 \end{aligned}$$

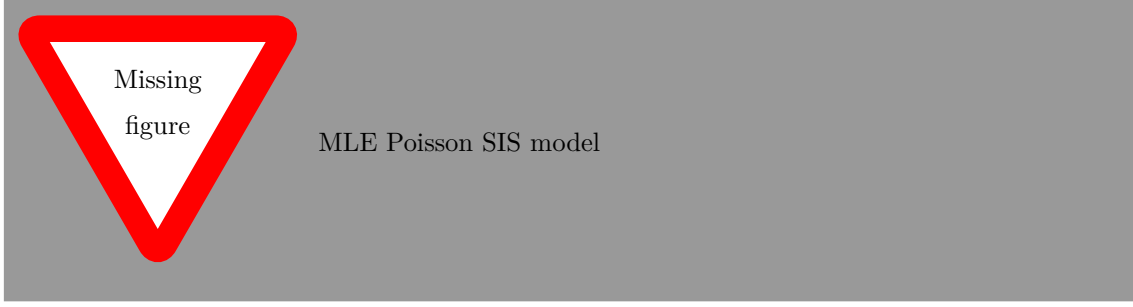
which we numerically solve to get $\hat{\theta}_0 \approx 1.329$ and $\hat{\theta}_1 \approx 0.751$, as seen in Figure 4.1.

Relationship of Least Squares and Maximum Likelihood Estimates

Although the least squares estimate does not explicitly assume a distribution, it coincides with the maximum likelihood estimate under the assumption that the y_i^{obs} s were has been generated with normal error.

Theorem 4.6. *If $f_i(\boldsymbol{\theta}) + \varepsilon_i \sim N(f_i(\boldsymbol{\theta}), \sigma^2)$, then*

$$\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}^{\text{LSE}}$$



Proof.

$$\begin{aligned}
 \hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta} \in \Theta} \ell(\boldsymbol{\theta}) \\
 &= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{(f_i(\boldsymbol{\theta}) - y_i^{\text{obs}})^2}{\sigma^2} \\
 &= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n -\frac{(f_i(\boldsymbol{\theta}) - y_i^{\text{obs}})^2}{\sigma^2} \\
 &= \arg \max_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n -(f_i(\boldsymbol{\theta}) - y_i^{\text{obs}})^2 \\
 &= \arg \min_{\boldsymbol{\theta} \in \Theta} \sum_{i=1}^n (f_i(\boldsymbol{\theta}) - y_i^{\text{obs}})^2 \\
 &= \boldsymbol{\theta}^{\text{LSE}}.
 \end{aligned}$$

□

Frequentist Parameter Estimates in Compartmental Models

Various approaches are possible to parameterise compartmental models. Often the ODE model is fit to the data directly (if there is only one datapoint) to fit to, or parameters are estimated by least squares. For example see Gani and Leach 2001 who use both these methods. including fitting the ODEs to by using the least squares estimates of a common approach is to find the least squares estimate of the parameters by considering $\mathbf{f}(\boldsymbol{\theta})$. This is commonly done because explicitly solving the likelihood for the stochastic compartmental model is infeasible. An alternative approach is to assume that the observed data follow a particular distribution determined by the ODE solution. For example, for the SIS model described in 2.1a it is plausible to assume that daily incidence (case counts) are distributed $\text{Pois}(\beta \frac{I_t}{N} S_t)$. Then β could be estimated using a maximum likelihood estimate given the case counts. Alternatively if case counts are overdispersed, a negative binomial likelihood could be used.

Chapter 5

Parameter Inference

TODO:

1. make all single digit numbers words

In classical statistics, Θ is fixed, and the data Y is generated from a distribution depending on Θ . Frequentist estimators such as the maximum likelihood estimator, or least squares estimator as previously seen are point estimates. In contrast, inference under a Bayesian framework assumes that Θ is also a random variable. This means although inference may be done on the same model, the Bayesian statistician will more often than not end up with a probability distribution describing Θ given Y .

5.1 Monte Carlo Integration

Let us assume we have a random variable $X \sim p$, and we want to integrate $I = \int_{\mathcal{X}} f(x)p(x) dx$ with relation to some known density p and function f . We can numerically approximate I by taking n i.i.d. samples $\{x_1, x_2, \dots, x_n\}$ from the distribution specified by p , and approximating I by $I_n := \frac{1}{n} \sum_{i=1}^n f(x_i)$. By the strong law of large numbers, under specific conditions (for example $\sigma_f := \text{Var}(f(X)) < \infty$ is sufficient) we have that $I_n \rightarrow I$ almost surely as $n \rightarrow \infty$. Furthermore by the central limit theorem we know that $\sqrt{n}(I_n - I) \sim N(0, \sigma_f)$ as $n \rightarrow \infty$. Intuitively, Monte Carlo integration works efficiently because the support set we numerically integrate over are likely to be points of higher probability density that integrating numerically over a set of uniformly distributed points.

5.2 Accept-Reject sampling methods

Accept-Reject methodology fundamentally relies on the theorem that $X \sim f$ is equivalent to simulating $(X, U) \sim \text{Unif}\{(x, u) | 0 < u < f(x)\}$. (Theorem 2.15 pg 47 Casella Roberts). In the case where a direct sampling method doesn't exist, we can use this theorem to our advantage. For example, if $a \leq X \leq b$ almost surely, then we can sample the pairs $X^* \sim \text{Unif}(a, b)$, and then $U^* \sim \text{Unif}(0, M)$, where $M \geq \sup_{a < x < b} f(x)$. To then get out samples to follow the distribution of (X, U) we then reject all samples (x_i, u_i) where $u_i > f(x_i)$. The distribution of X can then be described as follows

$$\begin{aligned}
\Pr(X \leq x) &= \Pr(X^* \leq x | U^* \leq f(X^*)) \\
&= \frac{\Pr(X^* \leq x, U^* \leq f(X^*))}{\Pr(U^* \leq f(X^*))} \\
&= \frac{\int_a^x f(y) dy / ((b-a) * M)}{1 / ((b-a) * M)} \\
&= \int_a^x f(y) dy \\
&= F(x)
\end{aligned}$$

For example in finding samples from the distribution below, we generate 100 random pairs (X^*, U^*)

There is no need for the distribution to be uniform over the rectangle, and to decrease the number of rejections we can generate points uniformly under a function $h(x)$ that more closely resembles our target distribution as long as $f(x) \leq h(x)$ for all x . If $h(x) = Mg(x)$ where $g(x)$ is a probability distribution function with known method of sampling, we can use this to generate via $Y \sim g$ and $U|Y \sim \text{Unif}(0, Mg(x))$ then we reject all pairs where $U|Y > f(Y)$.

So in the end we do

1. **Generate** $X \sim g$
2. **Accept** if $U \leq f(X)/Mg(X)$, **otherwise repeat** 1.

If f is computationally expensive, but bounded below by g_l , we can

1. **Generate** $X \sim g$
2. **Accept** if $U \leq g_l(X)/Mg(X)$
3. **Accept** if $U \leq f(X)/Mg(X)$, **otherwise repeat** 1.

5.3 Essentials for MCMC

Definition 5.1 (Random Walk). Random walk (MCSM 6.1) is a sequence $\{X_n\}$ such that

$$X_{n+1} = X_n + \varepsilon_n,$$

and ε_n is independent from X_1, \dots, X_n .

Definition 5.2 (Transition Kernel). A transition kernel is a function K defined on $\mathcal{X} \times \mathcal{B}(\mathcal{X})$ such that

1. $\forall x \in \mathcal{X}$, $K(x, \cdot)$ is a probability measure;
2. $\forall A \in \mathcal{B}(\mathcal{X})$, $K(\cdot, A)$ is measurable.

(MCSM 6.2)

Definition 5.3 (Markov Chain). Given a transition kernel K , a sequence of $X_0, X_1, \dots, X_n, \dots$ of random variables is a Markov chain denoted by (X_n) if for any t ,

$$P(X_{k+1} \in A | X_0, X_1, \dots, X_k) = P(X_{k+1} \in A | X_0, \dots, X_k)$$

(MCMS 6.4)

Definition 5.4 (φ -irreducible). *Given a measure φ , the Markov chain (X_n) with transition kernel K is φ -irreducible if, for every $B \in \mathcal{B}(\mathcal{X})$ with $\varphi(A) > 0$, there exists n s.t. $K^n(x, A) > 0$ for all $x \in \mathcal{X}$. (MCSM 6.13)*

Definition 5.5 (Atom). *The Markov chain (X_n) has an atom $\alpha \in \mathcal{B}(\mathcal{X})$ if there exists an associated nonzero measure ν such that*

$$K(x, A) = \nu(A)$$

(MCSM 6.18)

Definition 5.6 (Harris Recurrent). *A set A is Harris recurrent if $P_x(\eta_A = \infty) = 1$ for all $x \in A$. The chain (X_n) is Harris recurrent if there exists a measure ψ such that (X_n) is ψ -irreducible and for every set A with $\psi(A) > 0$, A is Harris recurrent. (MCSM 6.32)*

Definition 5.7 (Invariant (Measure)). *A σ -finite measure π is invariant for the transition kernel $K(\cdot, \cdot)$ if*

$$\pi(B) = \int_{\mathcal{X}} K(x, B) \pi(dx), \forall B \in \mathcal{B}(\mathcal{X})$$

(MCSM 6.35) *and π is a stationary distribution when π is a probability measure. If such a π exists for a ψ -reducible chain, the chain is positive.*

Definition 5.8 (Ergodic). *For a positive Harris recurrent chain (X_n) , with invariant distribution π , an atom α is ergodic if*

$$\lim_{n \rightarrow \infty} |K^n(\alpha, \alpha) - \pi(\alpha)| = 0$$

(MCSM 6.47)

5.4 Motivating the MH/Gibbs Algorithm

use some of the MC theory to arrive at the acceptance probability in MH

5.5 Metropolis Hastings

Definition 5.9 (Markov Chain Monte Carlo Method). *A Markov chain Monte Carlo method for simulating a distribution f is any method producing an ergodic Markov chain whose stationary distribution is f . (MCSM 7.1)*

Given a proposal distribution q , and a function f proportional to the posterior distribution, the Metropolis Hastings algorithm is as follows:

1. **Initialise** X_0
2. **Generate** $Y \sim q(y|x)$
3. **Calculate** $\alpha = \min\left(\frac{f(y)q(x|y)}{f(x)q(y|x)}, 1\right)$
4. $X_{t+1} = \begin{cases} Y, & \text{with probability } \alpha \\ X_t, & \text{else} \end{cases}$

If q is symmetric (that is $q(y|x) = q(x|y)$ for all x, y), then α simplifies down to $\frac{f(y)}{f(x)}$, and the algorithm is simply called the Metropolis algorithm.

As an example we can consider a series of coin tosses from a weighted coin. Let p be the probability of a head. We assume that $p \sim \text{Unif}(0, 1)$. After 10 tosses six heads could be observed. Let $H \sim \text{Binom}(10, p)$ be the number of heads that are tossed.

$$\begin{aligned}
\Pr(p \in dp | H = 6) &\propto \Pr(H = 6 | p \in dp) \Pr(p \in dp) \\
&= \binom{10}{6} p^6 (1-p)^{10-6} \\
&\propto p^6 (1-p)^4
\end{aligned}$$

Consider two different (symmetric) proposal distributions $q(p_i | p_{i-1})$:

- $p_i^* \sim N(p_{i-1}, 1/12)$
- $p_i^* \sim \text{Unif}(p_{i-1} - 1/2, p_{i-1} + 1/2)$.

Both are symmetric, and so

$$\begin{aligned}
\alpha &= \min \left(\frac{f(p_i) q(p_{i-1} | p_i)}{f(p_{i-1}) q(p_i | p_{i-1})}, 1 \right) \\
&= \min \left(\frac{p_i^6 (1-p_i)^4}{p_{i-1}^6 (1-p_{i-1})^4} \mathbb{I}(p_i > 0), 1 \right) \\
&= \min \left(\left(\frac{p_i}{p_{i-1}} \right)^6 \left(\frac{1-p_i}{1-p_{i-1}} \right)^4 \mathbb{I}(p_i > 0), 1 \right)
\end{aligned}$$

The effect of the choice of proposal distribution is shown in figure 5.1.

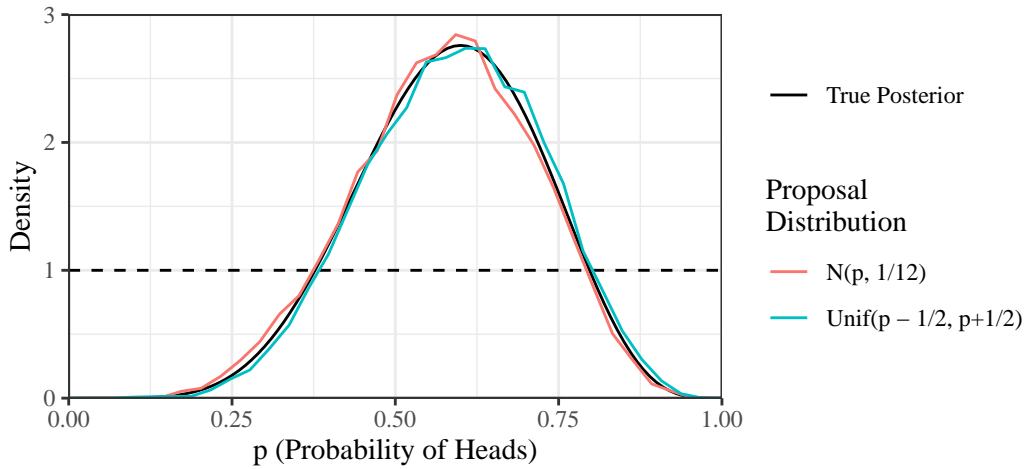


Figure 5.1: 30,000 samples from the posterior distribution of p using the Metropolis Hastings algorithm. It was assumed that $p \sim U(0, 1)$ and $H \sim \text{Binom}(10, p)$, given $H = 6$. A uniform and normal proposal distributions were compared.

For another example we can consider a simple SIS epidemiological model, with

$$\frac{dS}{dt} = \gamma I - \beta SI \quad \frac{dI}{dt} = \beta SI - \gamma I,$$

and a population of 25,000. Previous pandemics can be used to fix $\beta = 0.5$. At day 0, it is assumed that there is 1 infected person. Given 1184 new cases on day 30, and $\gamma \sim \text{Gamma}(2, 6)$ we can use the Metropolis Hastings algorithm to sample from the posterior distribution of γ . This will depend on the choice of likelihood. [Expand this explanation](#) See figure 5.2

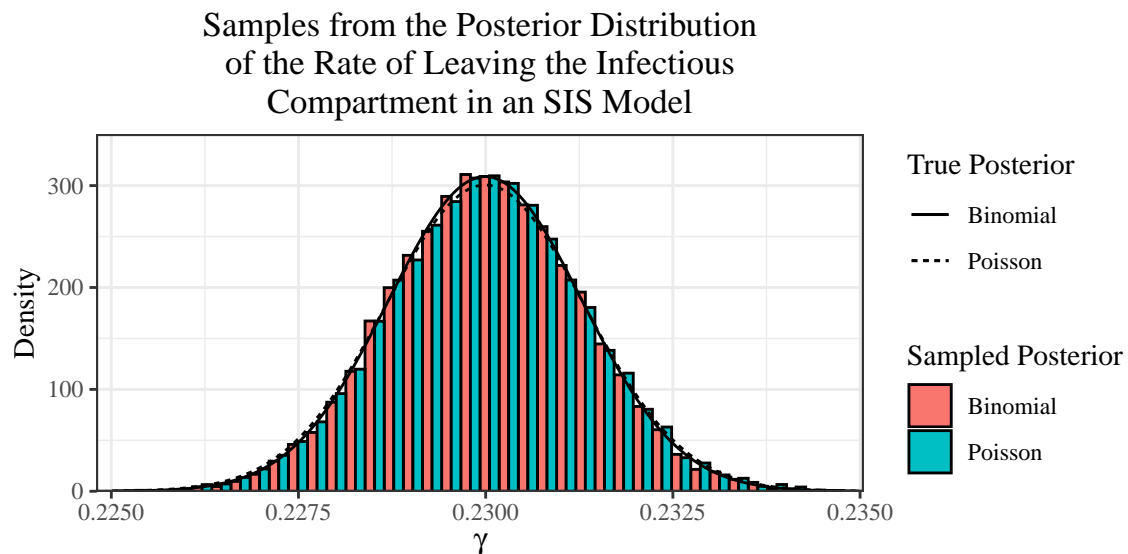


Figure 5.2: Using a basic SIS model, using incident changing likelihood function.

5.6 Gibbs's Sampling

For a multivariate case ($\Theta = \{\theta_1, \theta_2, \dots, \theta_n\}$), it can be analytically or computationally intractable to calculate the posterior (or even something proportional to the posterior) distribution. In this case it may be possible to calculate $\Pr(\theta_i \in d\theta | \theta_1, \theta_2, \dots, \theta_{i-1}, \theta_{i+1}, \theta_{i+2}, \dots, \theta_n, \mathbf{x})$ where \mathbf{x} is the data we are fitting the parameters to. **Prove Gibbs produces sample from the posterior distribution if this is in my final paper**

In this case we can use Gibbs sampling to sample from our posterior distribution. The basic algorithm is as follows:

Initialise for $\Theta = \{\theta_1, \theta_2, \dots, \theta_n\}$

1. For $i \in \{1, 2, \dots, n\}$ sample $\theta_i^* \sim \theta_i | \theta_1^*, \theta_2^*, \dots, \theta_{i-1}^*, \theta_{i+1}, \dots, \theta_n, \mathbf{x}$
2. Let $\Theta^* = \{\theta_1^*, \theta_2^*, \dots, \theta_n^*\}$
3. Append Θ^* to the chain of parameters, and let $\Theta := \Theta^*$

5.7 Diagnostics for Metropolis Hastings

include things like autoregression, thinning, burn-in etc.

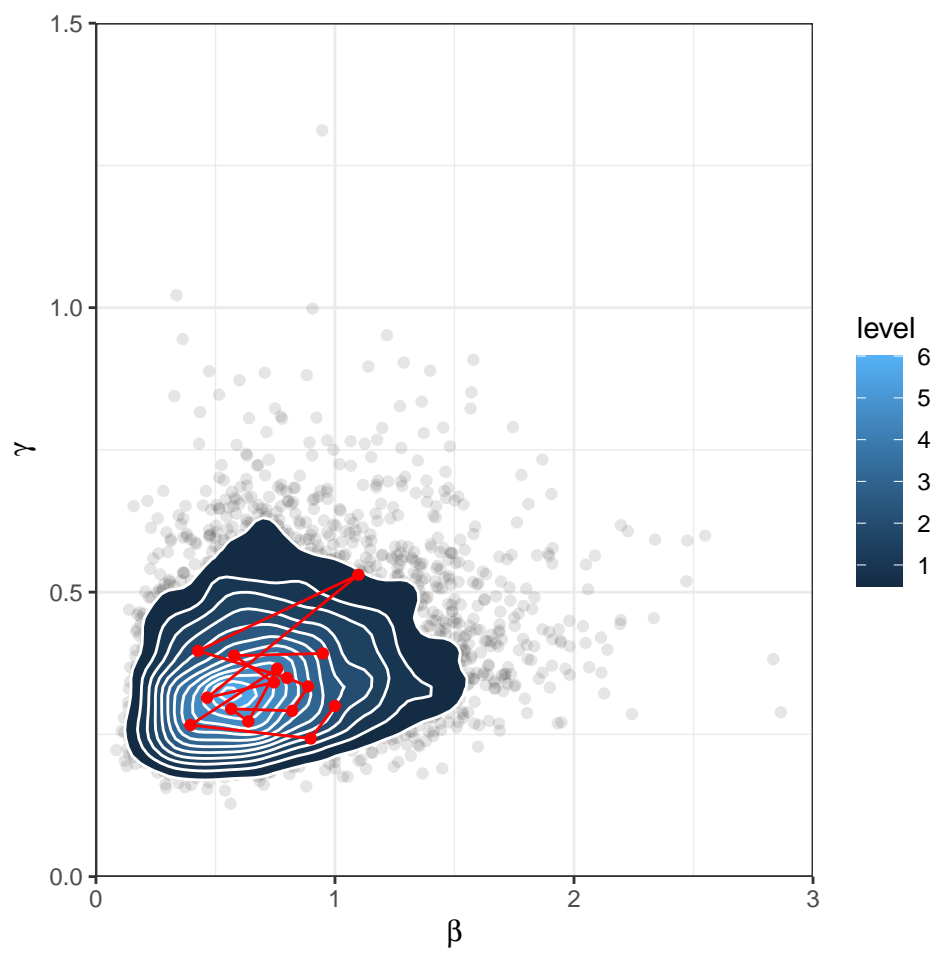


Figure 5.3: Using Gibbs on beta and gamma given an ' R_0 ' observation

Chapter 6

Gaussian Processes

- Prove all pos-sem is cov
- prove exponential quad kern is pos semi-def

6.1 Motivation and Definitions

If we knew the distribution of $D(\boldsymbol{\theta})$ exactly, then in order to do approximate Bayesian computation, running the model becomes superfluous to obtain a ‘sample’ from $D(\boldsymbol{\theta})$. Instead a sample $D(\boldsymbol{\theta})$ could be drawn directly from it’s distribution. It is highly improbable that the distribution of $D(\boldsymbol{\theta})$ is known in practice, and so this chapter describes a method of approximating the distribution. Furthermore, since $\Pr(D(\boldsymbol{\theta}) < \varepsilon)$ is approximately proportional to the true likelihood, sampling from the approximation of $D(\boldsymbol{\theta})$ can be used to for more efficient approximation of the likelihood. The approximation considered is achieved by modelling $D(\boldsymbol{\theta})$ as a realisation of a Gaussian process.

Definition 6.1 (Gaussian Process). *A collection of random variables $\{f(x)\}_{x \in \mathcal{X}}$ (where x may be a vector) is a Gaussian process if any finite subset of the collection of random variables is multivariate normal distributed. That is, there is a function $m : \mathcal{X} \rightarrow \mathbb{R}$ and symmetric kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that for all finite sets $\mathbf{x} := \{x_1, x_2, \dots, x_n\} \subset \mathcal{J}$, with $f(\mathbf{x}) := [f(x_1), f(x_2), \dots, f(x_n)]^T$*

$$f(\mathbf{x}) \sim \text{MVN} \left(\begin{bmatrix} m(x_1) \\ m(x_2) \\ \vdots \\ m(x_n) \end{bmatrix}, \mathbf{K} = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ k(x_n, x_1) & \dots & \dots & k(x_n, x_n) \end{bmatrix} \right).$$

Definition 6.2 (Mean and Covariance Function). *The the mean function and covariance kernel are*

$$m(x_i) := \mathbb{E}[f(x_i)]$$

and

$$k(x_i, x_{i'}) := \text{cov}(f(x_i), f(x_{i'})).$$

Although Gaussian processes are simultaneously realised over the whole space \mathcal{X} (for example \mathbb{R}^d) and are hence collections of (uncountably infinite) random variables, the choice of covariance

function $\text{corr}(x, x') \rightarrow 1$ as $\|x - x'\| \rightarrow 0$ induces continuity in x almost surely. Therefore they can be thought of as realisations of continuous functions. (Should I prove this?)

Some common examples of Gaussian processes include

1. Brownian motion on \mathbb{R} :

$$m \equiv 0, \quad \text{and} \quad k(s, t) = \min(s, t)$$

2. Ornstein Uhlenbeck process with parameters θ and σ :

$$m \equiv 0, \quad \text{and} \quad k(s, t) = \frac{\sigma_k^2}{2\theta} \left(e^{-\theta|t-s|} - e^{-\theta(t+s)} \right)$$

Properties such as the smoothness and undulation of the realised functions are also determined by the covariance kernel k , and associated hyperparameters. Before exploring different kernel options, we begin by defining a valid kernel function, and formalising ‘smoothness.’

Definition 6.3 (Positive Semi-Definite Matrix). *An $n \times n$ matrix \mathbf{A} is positive semi-definite if $\mathbf{v}^T \mathbf{A} \mathbf{v} \geq 0$ for all $\mathbf{v} \in \mathbb{R}^n$.*

Theorem 6.4 (Sufficient Condition for Positive Semi-Definite). *A symmetric matrix \mathbf{A} is positive semi-definite, if (and only if) it’s eigenvalues are non-negative.*

Proof. □

Definition 6.5 (Positive Semi-Definite Kernel). *A kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is positive semi-definite if the matrix*

$$\mathbf{K} = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ k(x_n, x_1) & \dots & \dots & k(x_n, x_n) \end{bmatrix}$$

is positive semi-definite for any collection of $x_i \in \mathcal{X}$

Theorem 6.6. *All symmetric positive semi-definite matrices are covariance matrices for some set of random variables*

Proof. □

Definition 6.7 (Mean Square Continuous). *A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is mean square continuous at \mathbf{x} in the i th direction at if $\mathbb{E}(|f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})|^2) \rightarrow 0$ as $|h| \rightarrow 0$, where \mathbf{e}_i is the unit vector with a 1 in the i th coordinate.*

Definition 6.8 (Mean Square Differentiable). *A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is mean square differentiable at \mathbf{x} in the i th direction with derivative $\frac{\partial f(\mathbf{x})}{\partial x_i}$ if*

$$\mathbb{E} \left[\left| \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h} - \frac{\partial f(\mathbf{x})}{\partial x_i} \right|^2 \right] \rightarrow 0$$

as $|h| \rightarrow 0$, where \mathbf{e}_i is the unit vector in the direction of the i th coordinate.

The concept of mean square differentiability and continuity are analogous to differentiability and continuity in the non-random function case.

Theorem 6.9. *Brownian motion is mean square continuous, but not mean square differentiable.*

Proof. $(B_{t+h} - B_t)^2 \sim (\sqrt{|h|}Z)^2$ where $Z \sim N(0, 1)$. Therefore $(B_{t+h} - B_t)^2 \sim |h|\chi_1^2 \rightarrow 0$ almost surely as $|h| \rightarrow 0$, and hence $\mathbb{E}[(B_{t+h} - B_t)^2] = 0$. Since $\frac{B_{t+h} - B_t}{h} \sim N(0, 1/|h|)$, $\frac{B_{t+h} - B_t}{h}$ does not converge to any valid probability distribution as $|h| \rightarrow 0$, as the variance approaches $+\infty$. \square

Some common

Theorem 6.10 (Positive Semi-Definiteness of RBF). *The radial basis function $\sigma_k^2 \exp(-\frac{(x-x')^2}{2\gamma^2})$ is a positive semi-definite kernel.*

Theorem 6.11 (Bochner's Theorem). *Let k be a stationary kernel function such that $k(x, x') = f(d)$. A function $k : \mathbb{R}^d \rightarrow \mathbb{C}$ is the covariance function of a weakly stationary mean square continuous complex-valued random process of \mathbb{R}^d if and only if it can be represented as*

$$k(\tau) = \int_{\mathbb{R}^d} \exp(2\pi i s \cdot \tau)$$

(Rasmussen and Williams 2008, p. 82)

6.2 Families of Kernel Function

The two most common families of kernel functions are the squared exponential and Matérn families.

Matérn Family

The Matérn exponential kernel is of the form

$$k_\nu(x, x') = \sigma_k^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - x'\|}{\ell} \right)^\nu K_\nu \left(-\frac{\sqrt{2\nu} \|x - x'\|}{\ell} \right)$$

where K_ν is a modified Bessel function (defined in Abramowitz and Stegun 2013, p. 374). The general form is not very insightful, however for $\nu = 1/2, 3/2$ and $5/2$, (the most common values used) the kernel can be written as:

$$\begin{aligned} k_{1/2}(x, x') &= \sigma_k^2 \exp \left(-\frac{\|x - x'\|}{\ell} \right) \\ k_{3/2}(x, x') &= \sigma_k^2 \left(1 + \frac{\sqrt{3}\|x - x'\|}{\ell} \right) \exp \left(-\frac{\sqrt{3}\|x - x'\|}{\ell} \right) \\ k_{5/2}(x, x') &= \sigma_k^2 \left(1 + \frac{\sqrt{5}\|x - x'\|}{\ell} + \frac{5\|x - x'\|^2}{3\ell^2} \right) \exp \left(-\frac{\|x - x'\|^2}{2 * \ell^2} \right) \end{aligned}$$

Zero mean Gaussian processes with a Matérn kernel are n times mean square differentiable, for all $n < \nu$. As seen in Figure 6.1, this means that this kernel allows for flexibility in how smooth realised functions are. As $\nu \rightarrow \infty$, with appropriate rescaling, the limit of the Matérn kernel is the squared exponential kernel. Rasmussen and Williams 2008, p. 85 **Proof in CHAPTER 4 SKOROKHOD STOCHASTIC I?**

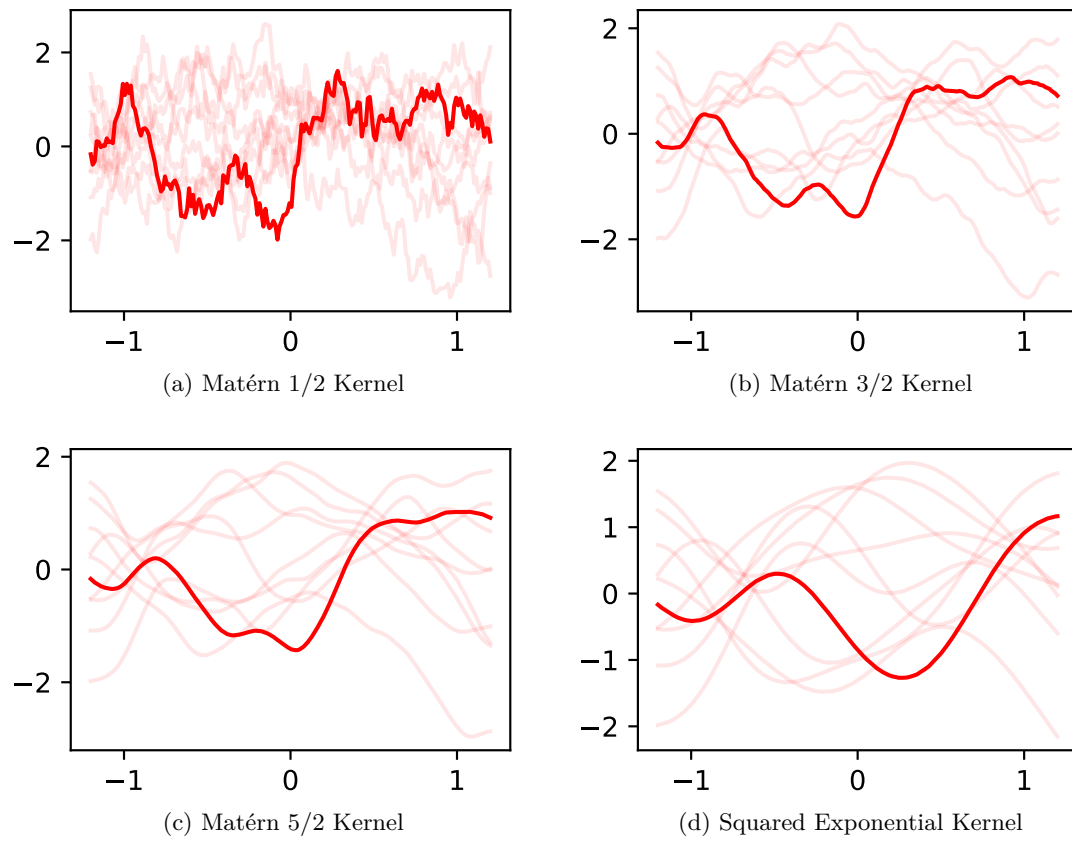


Figure 6.1: Ten sample realisations from 4 different kernels with hyperparameters $\ell = 1$, and $\sigma_o^2 = 1$. One realisation is bolded. Samples for each kernel were generated from the same seed.

Squared Exponential Kernel

The squared exponential kernel has the form

$$k(x, x') = \sigma_k^2 \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right)$$

As the limit of Matérn kernels, the squared exponential kernel is infinitely mean square differentiable. Despite this being the ‘default’ kernel in much of the literature, infinite differentiability is a very strong condition on functions which are very smooth, which can be seen in Figure 6.1d

Length and Amplitude Hyperparameters

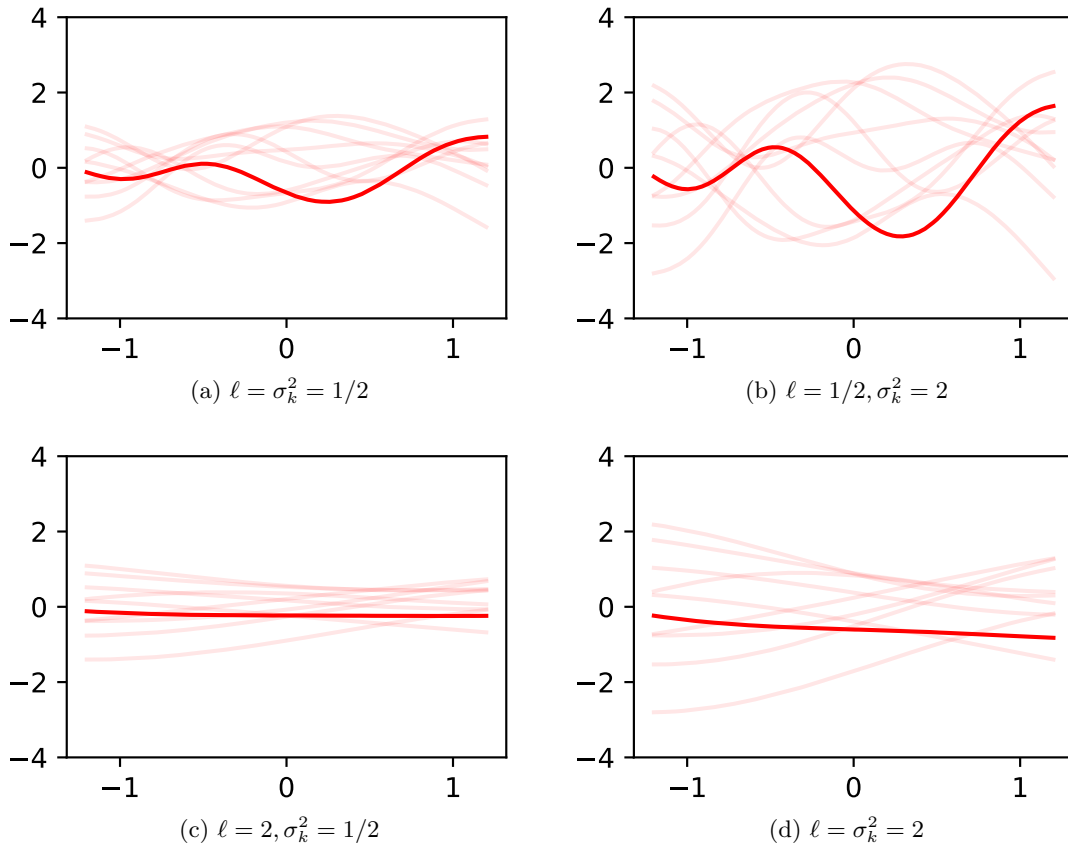


Figure 6.2: Ten realisations of zero mean Gaussian processes with the squared exponential kernel, varying the length and amplitude parameters. The samples were generated using the same seed

In both the Matérn and squared quadratic kernels (as well as most other common kernels choices), there are two hyperparameters ℓ and σ_k^2 which are referred to as length and amplitude hyperparameters. ℓ determines how close two points need to be to be highly correlated. Larger values of ℓ generates functions with higher correlation within a larger neighbourhood, as seen in Figure 6.2. σ_k^2 does not impact the correlation between x and x' , but scales the correlation matrix. In other words, larger σ_k^2 increase the size but not rate of fluctuations. This can be seen comparing Figure 6.2a to Figure 6.2b.

Chapter 7

Gaussian Process Regression

Given the set of observations $f(\mathbf{x}_*)$ for the set of indices \mathbf{x}_* , it is often desirable to infer information about the function values at unobserved values. By choosing a Gaussian process with fixed kernel and hyperparameters, we can condition the process on the observed data, limiting the family of possible functions that we assume truly describe the model. Under the assumption that the function is a realisation of a Gaussian process predicting unseen function values reduces to elementary linear algebra. This is because a conditional multivariate normal distribution is still multivariate normal, and so the distribution of unobserved points will be multivariate normal.

Consider

$$\begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}_*) \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m(\mathbf{x}) \\ m(\mathbf{x}_*) \end{bmatrix}, \begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix} \right)$$

Theorem 7.1 (Conditional Multivariate Normal Distribution is Multivariate Normal). *With $f(\mathbf{x})$ and $f(\mathbf{x}_*)$, the conditional distribution is*

$$f(\mathbf{x})|f(\mathbf{x}_*) \sim \mathcal{N} \left(m(\mathbf{x}) + K_* K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)), K - K_* K_{**}^{-1} K_*^T \right).$$

Proof. Since marginal distribution of the multivariate normal distribution, is also multivariate normal, $f(\mathbf{x}_*) \sim \mathcal{N}(m(\mathbf{x}_*), K)$. Let the inverse of $\begin{bmatrix} K & K_* \\ K_*^T & K_{**} \end{bmatrix}$ be defined as

$$\begin{bmatrix} \tilde{K} & \tilde{K}_* \\ \tilde{K}_*^T & \tilde{K}_{**} \end{bmatrix} = \begin{bmatrix} (K - K_* K_{**}^{-1} K_*^T)^{-1} & -(K - K_* K_{**}^{-1} K_*^T)^{-1} K_* K_{**}^{-1} \\ -K_{**}^{-1} K_*^T (K - K_* K_{**}^{-1} K_*^T)^{-1} & K_{**}^{-1} + K_{**}^{-1} K_*^T (K - K_* K_{**}^{-1} K_*^T)^{-1} K_* K_{**}^{-1} \end{bmatrix}$$

by the inverse of a block matrix. Therefore

$$\begin{aligned}
p(f(\mathbf{x})|f(\mathbf{x}_*)) &= \frac{p(f(\mathbf{x}), f(\mathbf{x}_*))}{p(f(\mathbf{x}_*))} \\
&\propto \frac{\exp \left[-\frac{1}{2} \left(\begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}_*) \end{bmatrix} - \begin{bmatrix} m(\mathbf{x}) \\ m(\mathbf{x}_*) \end{bmatrix} \right)^T \begin{bmatrix} K & K_* \\ K_*^T & K \end{bmatrix}^{-1} \left(\begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}_*) \end{bmatrix} - \begin{bmatrix} m(\mathbf{x}) \\ m(\mathbf{x}_*) \end{bmatrix} \right) \right]}{\exp \left[-\frac{1}{2} (f(\mathbf{x}_*) - m(\mathbf{x}_*))^T K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \right]} \\
&= \exp \left[-\frac{1}{2} \left(\begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}_*) \end{bmatrix} - \begin{bmatrix} m(\mathbf{x}) \\ m(\mathbf{x}_*) \end{bmatrix} \right)^T \begin{bmatrix} K & K_* \\ K_*^T & K \end{bmatrix}^{-1} \left(\begin{bmatrix} f(\mathbf{x}) \\ f(\mathbf{x}_*) \end{bmatrix} - \begin{bmatrix} m(\mathbf{x}) \\ m(\mathbf{x}_*) \end{bmatrix} \right) \right. \\
&\quad \left. + \frac{1}{2} (f(\mathbf{x}_*) - m(\mathbf{x}_*))^T K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \right] \\
&= \exp \left[-\frac{1}{2} \left((f(\mathbf{x}) - m(\mathbf{x}))^T \tilde{K} (f(\mathbf{x}) - m(\mathbf{x})) \right. \right. \\
&\quad \left. \left. + 2(f(\mathbf{x}) - m(\mathbf{x}))^T \tilde{K}_* (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \right. \right. \\
&\quad \left. \left. + (f(\mathbf{x}_*) - m(\mathbf{x}_*))^T \tilde{K}_{**} (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \right) \right. \\
&\quad \left. + \frac{1}{2} (f(\mathbf{x}_*) - m(\mathbf{x}_*))^T K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \right] \\
&\propto \exp \left[-\frac{1}{2} (f(\mathbf{x}) - m(\mathbf{x}))^T \tilde{K} (f(\mathbf{x}) - m(\mathbf{x})) \right. \\
&\quad \left. - (f(\mathbf{x}) - m(\mathbf{x}))^T \tilde{K}_* (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \right].
\end{aligned}$$

(by removing the terms independent of $f(\mathbf{x})$)

Since

$$p(\mathbf{z}) \propto \exp \left(-\frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z} + \mathbf{z}^T \mathbf{c} \right) \implies \mathbf{z} \sim \mathcal{N}(\Sigma \mathbf{c}, \Sigma),$$

$f(\mathbf{x}) - m(\mathbf{x})|f(\mathbf{x}_*)$ is multivariate normal with mean

$$\begin{aligned}
-\tilde{K}^{-1} \tilde{K}_* (f(\mathbf{x}_*) - m(\mathbf{x}_*)) &= (K - K_* K_{**}^{-1} K_*^T) \\
&\quad \times (K - K_* K_{**}^{-1} K_*^T)^{-1} K_* K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)) \\
&= K_* K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*))
\end{aligned}$$

and covariance matrix

$$\tilde{K}^{-1} = K - K_* K_{**}^{-1} K_*^T$$

by the alternative parametrisation of the multivariate normal distribution as a member of the exponential family of distributions (see Wikipedia contributors 2024, Table of Distributions). Finally, by the linearity of the multivariate normal mean,

$$f(\mathbf{x})|f(\mathbf{x}_*) \sim \mathcal{N} \left(m(\mathbf{x}) + K_* K_{**}^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)), K - K_* K_{**}^{-1} K_*^T \right).$$

□

After observing the function at multiple indices, we update the predictive distribution of any unobserved points, and generate new paths. The more points that the Gaussian process is conditioned on, the more narrow the sample paths, as seen in Figure 7.1

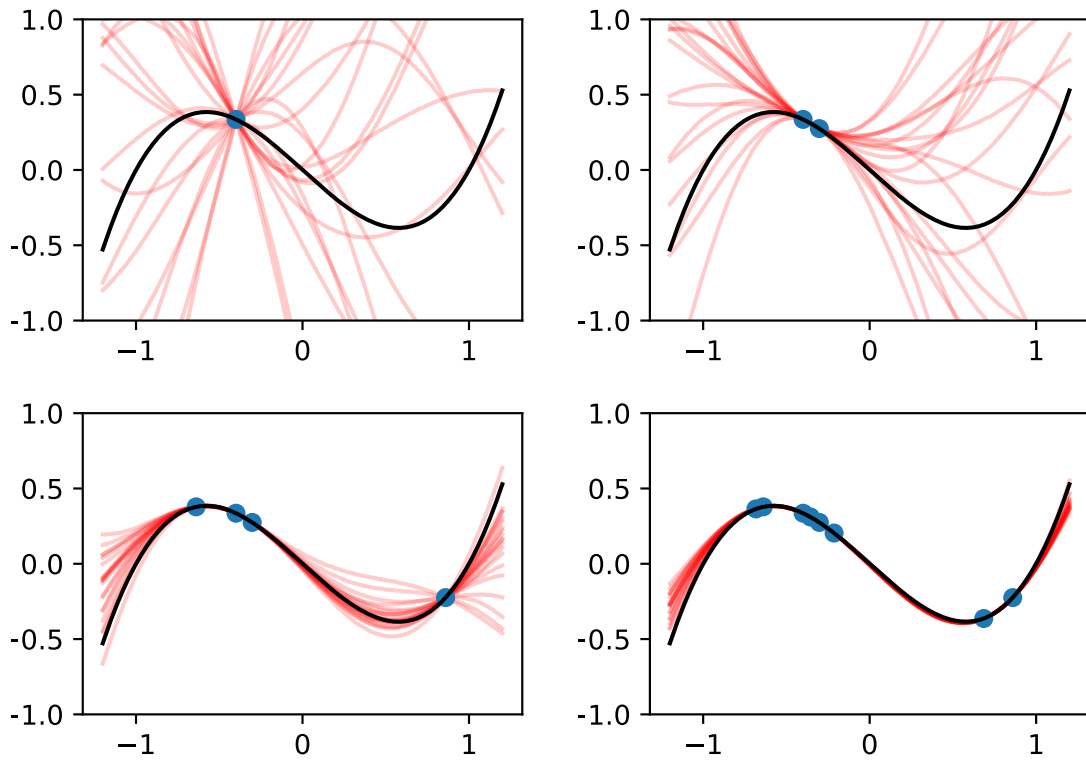


Figure 7.1: Sequence of Gaussian process regressions on the target function (black) $f(x) = x(x - 1)(x + 1)$, after 1, 2, 4, and 8 observations in blue. The red lines are new realisations from the conditioned Gaussian process. The Gaussian process was zero mean and had a squared exponential kernel. The hyperparameters were fixed with $\ell = 2.7$ and $\sigma_k^2 = 1.1$

7.1 Observation Variance

For most functions, model outputs, or processes desirable for approximating through Gaussian process regression, multiple observations (through model runs or an real life measurements) of the same point will result in different observations. This is in contrast to exact realisations as in Figure 7.1. The simplest assumption is that the observations are of the form

$$f_o(\mathbf{x}_*) = f(\mathbf{x}_*) + \varepsilon$$

where $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma_o^2 I)$. Under these assumptions, $\text{Cov}(f_o(\mathbf{x}_*), f_o(\mathbf{x}_*)) = K_{**} + \sigma_o^2 I$, where $K_{**} = \text{Cov}(f(\mathbf{x}_*), f(\mathbf{x}_*))$ matrix of $f(\mathbf{x}_*)$ without noise. Therefore the conditional distribution of our unobserved function outputs given noisy observations

$$f(\mathbf{x})|f_o(\mathbf{x}_*) \sim \mathcal{N}(m(\mathbf{x}) + K_*(K_{**} + \sigma_o^2 I)^{-1}(f(\mathbf{x}_*) - m(\mathbf{x}_*)), K - K_*(K_{**} + \sigma_o^2 I)^{-1}K_*^T).$$

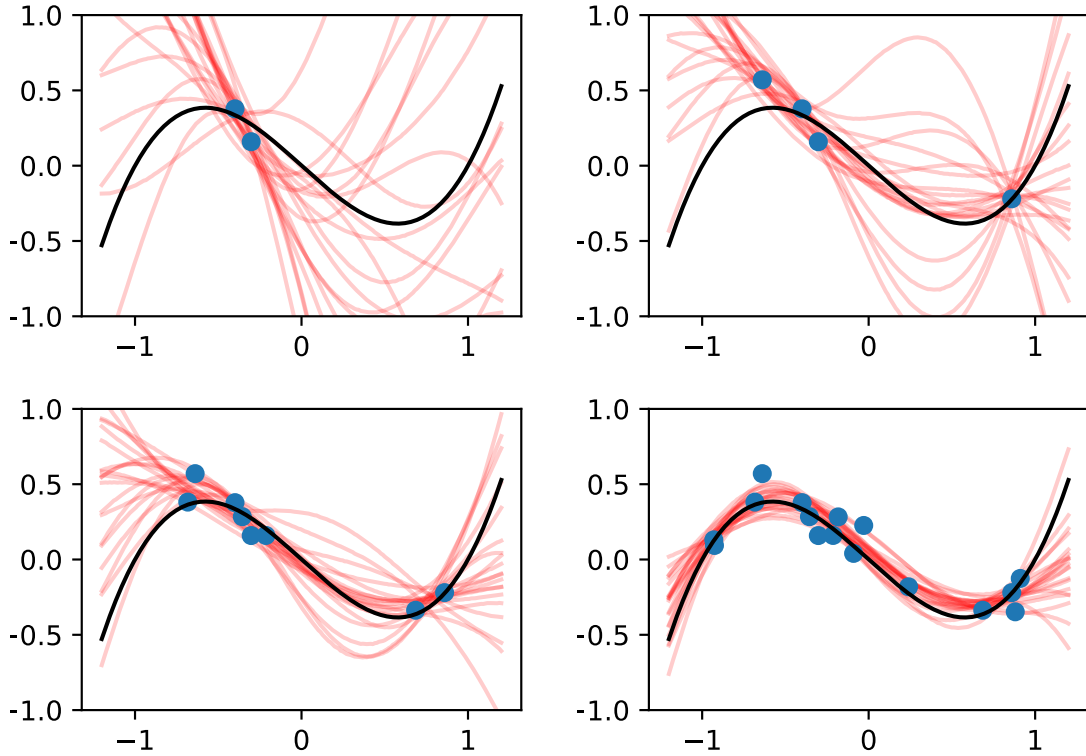


Figure 7.2: Sequence of Gaussian process regressions on the target function (black) $f(x) = x(x - 1)(x + 1)$, after 2, 4, 8, and 16 observations of $f(x_i) + \varepsilon_i$, where ε_i is i.i.d. $\mathcal{N}(0, \sigma_o^2)$ with $\sigma_o^2 = 0.01$ in blue. The red lines are new realisations from the conditioned Gaussian process. The Gaussian process was 0 mean and had a squared exponential kernel. The hyperparameters were fixed with $\ell = 2.7$ and $\sigma_k^2 = 1.1$

Adding noise to the observations makes the predictive distributions means that a single observation gives less information about the underlying function, and hence the predictive distributions of unseen data are much less stable, as empirically seen in Figure 7.2.

7.2 Model Selection

Kernel Family

The appropriate choice of kernel will depend on the properties behaviour of the target function to approximate. In the case of estimating an extremely stochastic distribution (such as the price of a stock over time), it is unlikely to be smooth, so no mean square differentiability is required, and a Matérn 1/2 kernel would be appropriate. If it is known that our target function is extremely smooth, such as a finite sum of infinitely differentiable functions (such as polynomials, sin, cos etc.) then the choice of squared exponential kernel is the most appropriate kernel. Realistically, the smoothness of the function will not be known a priori, and hence some sort of compromise (such as Matérn 5/2) kernel allows for flexibility.

Many other kernels exist that induce varying behaviours, such as periodic kernels [find periodic kernel](#), and non-stationary kernels (where the covariance is dependent on x and x' , not just $|x - x'|$) [have i defined stationary?](#).

Hyperparameters

The hyperparameters ℓ and σ_k^2 for a choice of kernel have to be are not known beforehand, unless the function is actually a realisation from a Gaussian process. Similarly, the observation variance σ_o^2 hyperparameter may not be a priori known. There are two main (frequentist) ways to fit these hyperparameters: maximum likelihood estimation, and leave-one-out cross validation.

Defining the likelihood $\mathcal{L}(\ell, \sigma_k^2, \sigma_o^2) := p(f(\mathbf{x}_*) | \ell, \sigma_k^2, \sigma_o^2)$ in the usual way, the maximum likelihood estimates are

$$\{\hat{\ell}, \hat{\sigma}_k^2, \hat{\sigma}_o^2\} := \arg \max_{\{\ell, \sigma_k^2, \sigma_o^2\}} \mathcal{L}(\ell, \sigma_k^2, \sigma_o^2)$$

which is equivalent to minimising

$$-\ln(\mathcal{L}) = \frac{1}{2} [\ln(|K_{**}(\ell, \sigma_k^2) + \sigma_o^2|) + (f(\mathbf{x}_*) - m(\mathbf{x}_*))^T (K_{**}(\ell, \sigma_k^2) + \sigma_o^2)^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*)) + c] .$$

The covariance matrix generated by the choice of kernel K_{**} is explicitly written with its dependence on ℓ and σ_k^2 . c is a constant.

Leave-one-out cross validation aims to maximise the predictive log probability.

$$\{\tilde{\ell}, \tilde{\sigma}_k^2, \tilde{\sigma}_o^2\} := \arg \max_{\ell, \sigma_k^2, \sigma_o^2} \sum_i \ln p(f_i(\mathbf{x}_*) | f_{-i}(\mathbf{x}_*), \ell, \sigma_k^2, \sigma_o^2),$$

where $f_i(\mathbf{x}_*) | f_{-i}(\mathbf{x}_*)$ is the distribution of the i th element of $f(\mathbf{x}_*)$ conditioned on the rest of the observed data excluding that element (represented by $f_{-i}(\mathbf{x}_*)$). $f_i(\mathbf{x}_*) | f_{-i}(\mathbf{x}_*)$ can be found by Theorem 7.1. Computationally efficient methods for calculating the predictive log probability that avoid having to invert the covariance matrix for every summand element exist. In particular it can be shown that $f_i(\mathbf{x}_*) | f_{-i}(\mathbf{x}_*)$ has mean

$$f_i(\mathbf{x}_*) - m_i(\mathbf{x}_*) - [(K_{**} + \sigma_o^2 I)^{-1} (f(\mathbf{x}_*) - m(\mathbf{x}_*))]_i / [(K_{**} + \sigma_o^2 I)^{-1}]_{ii}$$

and variance $1 / [(K_{**} + \sigma_o^2 I)^{-1}]_{ii}$, where both the mean and covariance are (surprisingly) independent of $f_i(\mathbf{x}_*)$ (Rasmussen and Williams 2008, p. 116).

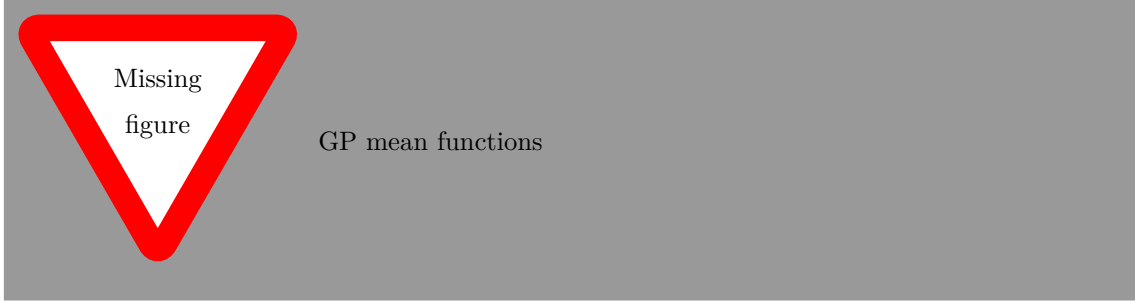
Both methods can be extended to include estimating hyperparameters given a family of possible

mean functions. For example Gutmann and Cor 2016 use maximum likelihood estimates for the amplitude and turning points of the quadratic mean functions.

Recent work has shown that at least under specific conditions, the leave-one-out estimates for the scale hyperparameter are more robust to a larger family of target functions (Naslidnyk et al. 2024), and the broader literature seems to favor leave-one-out cross validation.

Finally there is scope for a Bayesian approach to model selection. By setting priors on the hyperparameters to be estimated and using the likelihood as described in the maximum likelihood estimation approach, a posterior distribution can be easily contrived. A set of samples from this posterior could then be taken, or a less principled maximum a posteriori probability estimate could then be taken for a point estimate of the hyperparameters. This approach has some obvious benefits, particularly when taking a posterior sample. Most obviously is that since multiple values for each hyperparameter are sampled, the set of functions after

7.3 Differing mean functions



7.4 Bayesian Acquisition Functions

Under the assumptions that making observations from the underlying function is costly, and we care about regions of the function with high (or low) values, new observations should be taken where there is high probability the function will be low. There also needs to be a trade-off between observing from areas with high predictive mean, and high predictive variance. These ideas are formalised by Bayesian acquisition functions $\mathcal{A}(x)$, with larger values corresponding to a higher ‘desirability.’ The target function is then sampled at the x which maximises this acquisition function. The new observation is then incorporated into the acquisition function.

Upper Confidence Bound

The upper confidence bound is one common way of exploring this trade off. The upper confidence bound

$$\mathcal{A}_{\text{UCB}}(x) := \mathbb{E}[f(x)|f(\mathbf{x}_*)] + \eta_t \sqrt{\text{Var}[f(x)|f(\mathbf{x}_*)]}$$

$$\mu(\boldsymbol{\theta}) - \eta_t \sqrt{\text{v}(\boldsymbol{\theta})}$$

with $\eta_t := \sqrt{2 \ln(\frac{t^{2d+2}\pi^2}{3\varepsilon})}$, and $\varepsilon \in (0, 1)$ that can be chosen (with a lower epsilon resulting in more exploration), $\mu(\boldsymbol{\theta})$ and $\text{v}(\boldsymbol{\theta})$ are the posterior mean and variance. $\varepsilon = 0.1$ was used.

Similarly we can also define

- BOLFI paper uses

$$\mu(\boldsymbol{\theta}) - \eta_t \sqrt{v(\boldsymbol{\theta})}$$

- $\eta_t := \sqrt{c + 2 \ln(t^{d/2+2})}$, and c can be chosen
- $\mu(\boldsymbol{\theta})$ and $v(\boldsymbol{\theta})$ are the posterior mean and variance

- Could use expected information

$$(\mu_{\min} - \mu(\boldsymbol{\theta}))\Phi\left(\frac{\mu_{\min} - \mu(\boldsymbol{\theta})}{\sqrt{v(\boldsymbol{\theta})}}\right) + \sqrt{v(\boldsymbol{\theta})}\phi\left(\frac{\mu_{\min} - \mu(\boldsymbol{\theta})}{\sqrt{v(\boldsymbol{\theta})}}\right)$$

- $\mu_{\min} := \min_{\boldsymbol{\theta}} \mu(\boldsymbol{\theta})$
- Φ, ϕ CDF and PDF of standard normal

exploration parameter proven by Srinivas et al. 2010

Show why Bayesian Acquisition fn not work for infinite space.

Part II

Calibrating Parameters for a *P.* *vivax* Model

Chapter 8

Methods

8.1 Creation of Synthetic Data

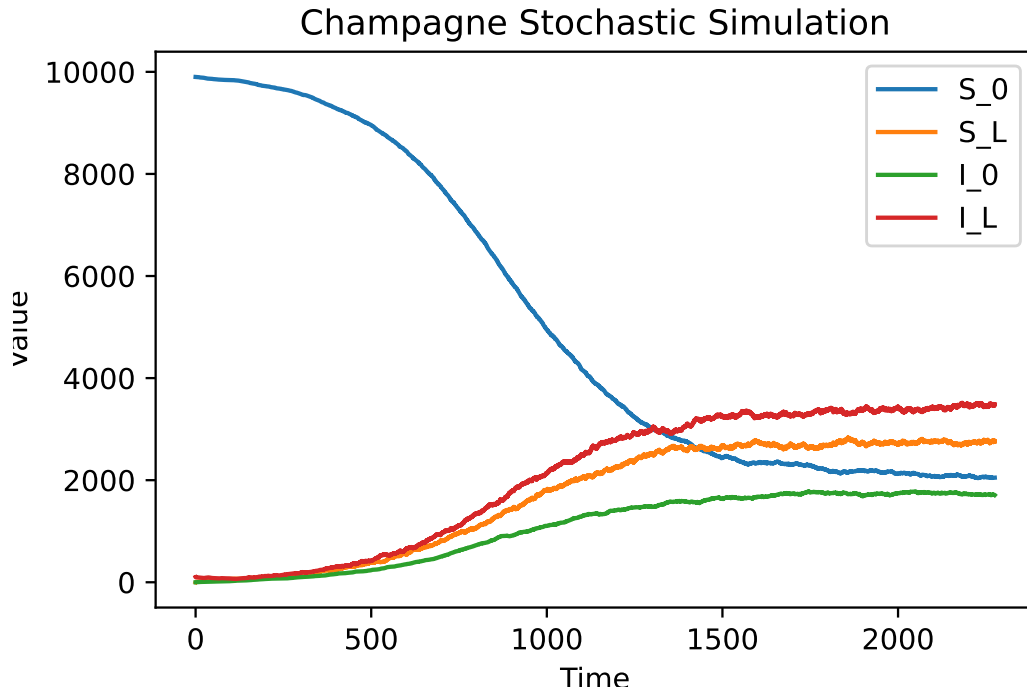


Figure 8.1: A Doob-Gillespie Simulation of the model described by Champagne et al. 2022 with $\alpha = 0.4$, $\beta = 0.4$, $\gamma_L = 1/223$, $\lambda = 0.04$, $f = 1/72$, $r = 1/60$, and $\delta = 0$. The population was 1000, with 10 initial infections (both blood and liver stage).

We investigated the model by Champagne et al. 2022 as described in 3.2. A malaria epidemic was simulated using the Doob-Gillespie algorithm (see Figure 8.1), using a population size of 1,000, and initial infected population of 10 (with both liver and blood stage infection). The parameters used closely followed those reported in Champagne et al. 2022:

- effective blood stage treatment proportion $\alpha = 0.4$,
- effective liver stage treatment proportion $\beta = 0.4$,

- rate of liver stage disease clearance $\gamma_L = 1/223 \text{ days}^{-1}$,
- importation rate $\delta = 0$ (assumed to be known),
- rate of infection $\lambda = 0.04 \text{ days}^{-1}$,
- rate of relapse $f = 1/72 \text{ days}^{-1}$,
- rate of blood stage disease clearance $r = 1/60 \text{ days}^{-1}$.

From intialisation, the simulation was run for 15,000 events (with an event being anything that caused the size of any compartment to change such as an infection, recovery, relapse etc.), after which, the model was understood to have reached steady state behaviour.

The synthetic data (as summary statistics) measured from the steady state were:

1. p_{obs} : the number of currently infected individuals at the end of simulated epidemic (steady state prevalence).
2. m_{obs} : the number of cases in the first week of the epidemic (first month incidence)
3. w_{obs} : the number of cases in the last week of the epidemic (steady state weekly incidence).

New infections which instantly undergo radical cure don't change the size of each compartment. The number of these 'silent' incidences were calculated between events using a Poisson distribution with rate $\Delta t \times \alpha\beta\lambda(I_L + I_0)S_0/N$, where Δt is the time between events.

8.2 Model Simulations and Discrepancy Function

New epidemics were simulated as above, with 15,000 events and at least 30 days (to allow for calculation of incidence in the first month of the epidemic). For each model run steady state prevalence p , first month incidence m , and steady state weekly incidence w were calculated with the identical method to the synthetic data.

We defined the discrepancy function to be L_2 norm of the relative differences

$$\mathcal{D}(\alpha, \beta, \gamma_L, \lambda, f, r) := \sqrt{\left(\frac{p - p_{\text{obs}}}{p_{\text{obs}}}\right)^2 + \left(\frac{m - m_{\text{obs}}}{m_{\text{obs}}}\right)^2 + \left(\frac{w - w_{\text{obs}}}{w_{\text{obs}}}\right)^2}.$$

Relative difference was chosen to limit the impact between the scale differences of the summary statistics.

8.3 Gaussian Process and Initialisation

We used a Gaussian process $d_{\mathcal{GP}}(\boldsymbol{\theta})$ as a surrogate model for $\mathbb{E}[\ln \mathcal{D}(\boldsymbol{\theta})]$. It was regressed on samples of $\overline{\ln \mathcal{D}(\boldsymbol{\theta})}$, where $\overline{\ln \mathcal{D}(\boldsymbol{\theta})}$ is the sample mean of 30 $\ln \mathcal{D}(\boldsymbol{\theta})$ samples. We used the kernel

$$k(\boldsymbol{\theta}_i, \boldsymbol{\theta}'_i) = \sigma_k^2(1 + z_i + \frac{z_i^2}{3})\exp(-z_i)$$

where

$$z_i = \sqrt{5 \sum_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \left(\frac{\theta_i - \theta'_i}{\ell_{\boldsymbol{\theta}}} \right)^2};$$

that is, a Matérn kernel with $\nu = 5/2$ and automatic relevance determination - i.e. each parameter $\theta \in \boldsymbol{\theta}$ was scaled by ℓ_θ . This can be seen as giving each parameter its own length hyperparameter. The Gaussian process was constant mean $m_{\mathcal{GP}}$ with ‘observation’ variance σ_o^2 . σ_o^2 is naturally interpretable as the sample mean variance.

Modelling the mean of the (log) discrepancy function is theoretically more sound than modelling the (log) discrepancy function itself, because we were not able to find theoretically sound reasons to assume normality, but the mean is asymptotically normal under reasonable conditions by the central limit theorem. For example, even under the simplest case that the discrepancy is L_2 norm of k differences x_i where $x_i \sim \mathcal{N}(0, 1)$, $\sqrt{\sum_{i=1}^k x_i^2} \sim \chi(k)$, a Chi distribution with k degrees of freedom. [do I need to cite a reference for this distribution \(wikipedia\)](#)

| Parameter | Upper Bound | Unit |
|---|-------------|--------|
| Proportion of treatment clearing blood stage disease α | 1 | |
| Proportion of treatment clearing liver stage disease β | 1 | |
| Rate of liver stage disease clearance γ_L | 1/30 | 1/days |
| Rate of infection λ | 1/10 | 1/days |
| Rate of relapse f | 1/14 | 1/days |
| Rate of blood stage disease clearance r | 1/14 | 1/days |

Table 8.1: Conservative upper bounds for parameters to be calibrated. Values were informed by Champagne et al. 2022; White et al. 2016. All lower bounds were zero.

All parameters to be calibrated were given conservative upper bounds after considering values reported in the literature, which informed where to define the Gaussian process. Parameter values outside this range were not considered.

Latin hypercube sampling was used to generate initialise 50 samples of the parameter space (scaled to be between zero and the upper bounds described in Table 8.1). For each set of parameters, $\overline{\ln \mathcal{D}(\boldsymbol{\theta})}$ was generated. The hyperparameters $m_{\mathcal{GP}}$, σ_o^2 , σ_a^2 , ℓ_α , ℓ_β , ℓ_{γ_L} , ℓ_λ , ℓ_f , and ℓ_r were selected by leave one out cross validation.

8.4 Bayesian Acquisition and Parameter Updates

For 500 iterations, the next $\boldsymbol{\theta}$ to sample $\overline{\ln \mathcal{D}(\boldsymbol{\theta})}$ from at was found by maximising the expected information acquisition function $\mathcal{A}_{\text{EI}}(\boldsymbol{\theta})$. Let the current iteration be t , and $d_{\mathcal{GP}}^{(i)}(\boldsymbol{\theta})$ be the Gaussian process regressed on the simulated $\overline{\ln \mathcal{D}(\boldsymbol{\theta})}$ s after i iterations.

Each index of $\boldsymbol{\theta}$ was initialised randomly at either the previous sample of $\boldsymbol{\theta}$ which minimised $\mathbb{E}(d_{\mathcal{GP}}^{(t-1)}(\boldsymbol{\theta}))$, or uniformly at random between it’s lower and upper bounds. In each iteration, there was a $\min[1/5 + \exp(1 - t/4), 1]$ probability that one of the parameters (say θ^*) in $\boldsymbol{\theta}$ was chosen, and $\overline{\ln \mathcal{D}(\boldsymbol{\theta})}$ was sampled at $\boldsymbol{\theta}$ as well as at 11 evenly spaced values of θ^* , with the other parameters fixed. $1/5 + \exp(1 - t/4) > 1$ for small t , decaying to $1/5$ as t is large. This was to help initialise the Gaussian process model, as well as optimise the ℓ_θ s. Every 50 iterations, the hyperparameters were reoptimised using leave one out cross validation. Finally, the synthetic likelihood was calculated using $\hat{\mathcal{L}}(\boldsymbol{\theta}) := \Pr(d_{\mathcal{N}}(\boldsymbol{\theta}) < \epsilon)$, where $d_{\mathcal{N}} \sim \mathcal{N}(\mathbb{E}[d_{\mathcal{GP}}^{(500)}(\boldsymbol{\theta})], 30^2 \sigma_o^2)$

The Gaussian process and Gaussian process regression was implemented using TensorFlow in Python (Martín Abadi et al. 2015).

Chapter 9

Results and Discussion

9.1 Results

should these be presented in tables?

The synthetic data generated were

- $p_{\text{obs}} = 593$ final prevalence
- $m_{\text{obs}} = 13$
- $w_{\text{obs}} = 73$.

The final hyperparameter estimates were

- $\sigma_k^2 = 0.677$
- $\sigma_o^2 = 0.093$
- $\ell_\alpha = 0.232$
- $\ell_\beta = 0.25$
- $\ell_{\gamma_L} = 0.008$
- $\ell_\lambda = 0.015$
- $\ell_f = 0.018$
- $\ell_r = 0.013$

The minimum predicted mean of a sampled θ was at (0.492, 0.9, 0.029, 0.03, 0.064, 0.023), compared to the true parameters used to generate the data (0.4, 0.4, 0.004, 0.04, 0.014, 0.017).

As the number of Bayesian acquisitions increased, the Gaussian process mean seems to converge to the true mean, as seen in Figure 9.1. The likelihood has maximum around the true values as seen in Figure 9.2. I will add in lines where the true values are to this figure. There seems to be a lack of good data

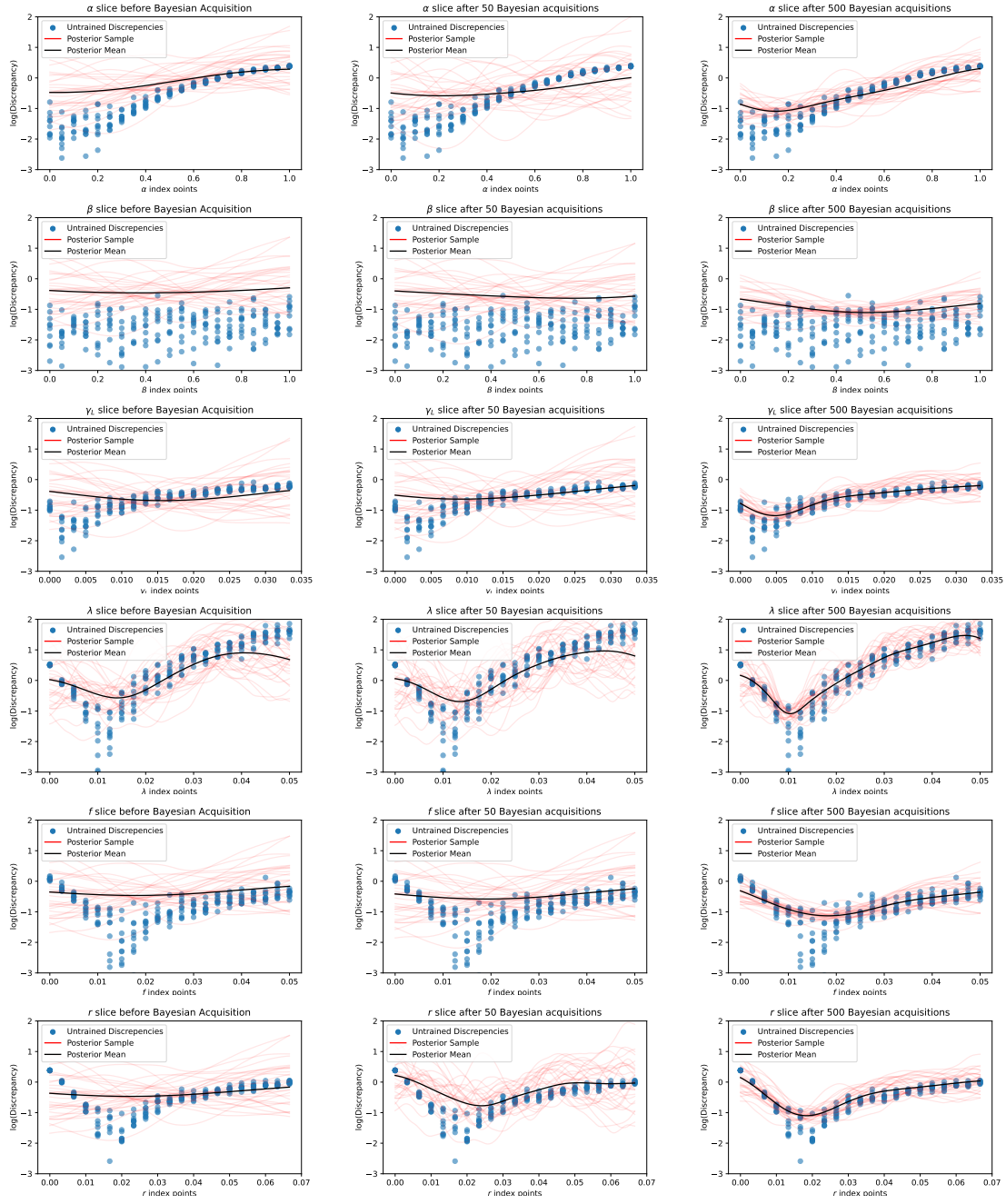


Figure 9.1: Four sample discrepancies $\ln \mathcal{D}(\theta)$ taken along 21 values of each parameter to be predicted between the lower and upper parameter bound. The Gaussian process predicts the mean of $\ln \mathcal{D}(\theta)$ after initialisation, 50 iterations of Bayesian acquisition, and 500 iterations of Bayesian acquisition. All parameters not in the slice are fixed at the true parameters that generated the synthetic observed data. The black lines are $\mathbb{E}(d_{\mathcal{GP}}(\theta))$ and the red lines are sample realisations from $(d_{\mathcal{GP}}(\theta))$

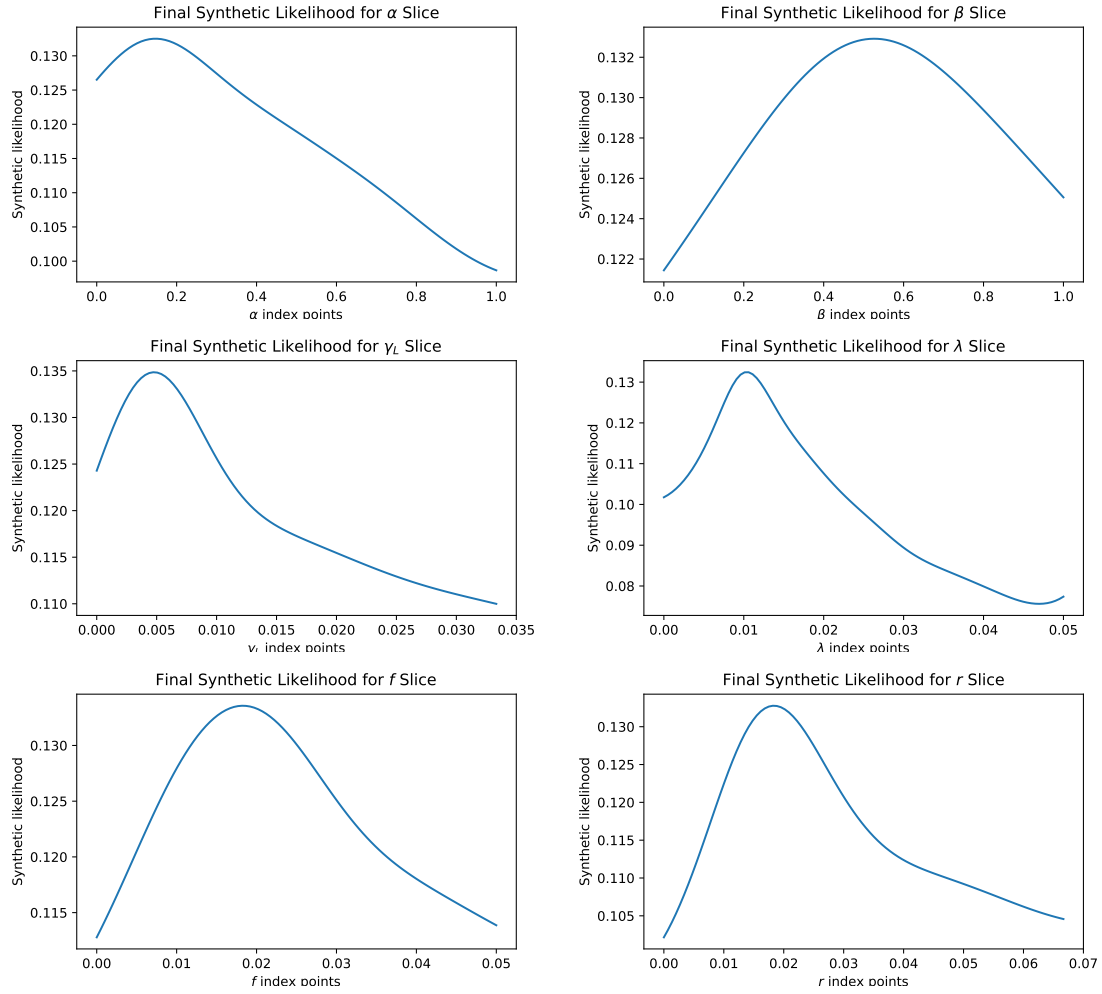


Figure 9.2: Final synthetic likelihoods $\hat{L}(\theta)$ after 500 iterations of Bayesian acquisition. Each likelihood shows how that likelihood changes across the parameter space.

9.2 Discussion

Although most of the procedure we used closely follows the paper by Gutmann and Cor 2016, there are a few key ways in which we modified the manuscript’s method. In particular, we felt the choice of squared exponential kernel was not a good assumption, since it implicitly assumes a high degree of smoothness in the target discrepancy function. This is unlikely to be met if the model has any bifurcation points. Under the squared exponential function, if the true function has any sharp declines, or general non-smoothness, the length scale is forced to be set very small. Therefore, although the squared exponential is the most commonly chosen kernel, we used a Matérn kernel with $\nu = 5/2$. This is not as constrained as a squared exponential kernel, but realisations are twice mean square differentiable.

Gutmann and Cor 2016 use the lower confidence bound acquisition function, where the exploration parameter is written as $\eta_t := \sqrt{2 \ln(\frac{t^{2d+2}\pi^2}{3\varepsilon})}$. This seems to be a mistake inherited from Brochu, Cora, and Freitas 2010. Even using Brochu’s original citation (Srinivas et al. 2010) to revise this to $\eta_t := \sqrt{2 \ln(\frac{t^{2d+2}\pi^2}{3\varepsilon})}$.¹ When we tried both of these exploration parameters, the choice of ε between $(0, 1)$ largely lead to repeated sampling from the same set of parameters, even for very small ε .

This is problematic, as Srinivas et al. 2010 assume a compact parameter space, whereas Gutmann and Cor 2016 do not.

They use a quadratic mean structure on unbounded parameters. This could be problematic in situations where the discrepancy has multiple local minima. When we tried it, the quadratic mean function poorly it resulted in in the acquisition function trying to sample extremely high parameter values. Therefore we bounded our parameters, making out parameter space compact.

9.3 Further Work

Moment matching could be used to avoid assuming that the distribution of the discrepancy $\mathcal{D}(\theta)$ or $\ln \mathcal{D}(\theta)$ is normal for fixed $\mathcal{D}(\theta)$. The sample variance could be explicitly modelled using

¹One Python package that implements BOLFI notes this error, see: <https://github.com/elfi-dev/elfi/blob/dev/elfi/methods/bo/acquisition.py>

Bibliography

- Abramowitz, Milton and Irene A. Stegun, eds. (2013). *Handbook of mathematical functions: with formulas, graphs, and mathematical tables*. 9. Dover print.; [Nachdr. der Ausg. von 1972]. Dover books on mathematics. New York, NY: Dover Publ. 1046 pp. ISBN: 978-0-486-61272-0.
- Adams, John H. and Ivo Mueller (Sept. 2017). “The Biology of Plasmodium vivax”. In: *Cold Spring Harbor Perspectives in Medicine* 7.9, a025585. ISSN: 2157-1422. DOI: 10.1101/cshperspect.a025585. URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5580510/> (visited on 03/24/2023).
- Aron, Joan L. and Robert M. May (1982). “The population dynamics of malaria”. In: *The Population Dynamics of Infectious Diseases: Theory and Applications*. Ed. by Roy M. Anderson. Boston, MA: Springer US, pp. 139–179. ISBN: 978-1-4899-2901-3. DOI: 10.1007/978-1-4899-2901-3_5. URL: https://doi.org/10.1007/978-1-4899-2901-3_5.
- Brochu, Eric, Vlad M. Cora, and Nando de Freitas (2010). *A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning*. arXiv: 1012.2599 [cs.LG].
- Champagne, Clara et al. (Jan. 2022). “Using observed incidence to calibrate the transmission level of a mathematical model for Plasmodium vivax dynamics including case management and importation”. In: *Mathematical Biosciences* 343, p. 108750. ISSN: 00255564. DOI: 10.1016/j.mbs.2021.108750. URL: <https://linkinghub.elsevier.com/retrieve/pii/S0025556421001541> (visited on 08/22/2023).
- Cowman, Alan F. et al. (2016). “Malaria: Biology and Disease”. In: *Cell* 167.3. Type: Review, pp. 610–624. DOI: 10.1016/j.cell.2016.07.055. URL: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-84994000411&doi=10.1016%2fj.cell.2016.07.055&partnerID=40&md5=81d9b4c51fe738ac66e0c8561b12c5bf>.
- Gani, Raymond and Steve Leach (Dec. 13, 2001). “Transmission potential of smallpox in contemporary populations”. In: *Nature* 414.6865, pp. 748–751. ISSN: 0028-0836, 1476-4687. DOI: 10.1038/414748a. URL: <https://www.nature.com/articles/414748a> (visited on 06/10/2024).
- Gutmann, Michael U. and Jukka Cor (2016). “Bayesian Optimization for Likelihood-Free Inference of Simulator-Based Statistical Models”. In: *Journal of Machine Learning Research* 17.125, pp. 1–47. ISSN: 1533-7928. URL: <http://jmlr.org/papers/v17/15-017.html> (visited on 04/28/2024).
- Hagenaars, T. J., C. A. Donnelly, and N. M. Ferguson (Apr. 2006). “Epidemiological analysis of data for scrapie in Great Britain”. en. In: *Epidemiology and Infection* 134.2, pp. 359–367. ISSN: 0950-2688, 1469-4409. DOI: 10.1017/S0950268805004966. URL: https://www.cambridge.org/core/product/identifier/S0950268805004966/type/journal_article (visited on 03/26/2024).

- Keeling, Matthew James and Pejman Rohani (2008). *Modeling infectious diseases in humans and animals*. OCLC: ocn163616681. Princeton: Princeton University Press. 366 pp. ISBN: 978-0-691-11617-4.
- Martín Abadi et al. (2015). *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*. Software available from tensorflow.org. URL: <https://www.tensorflow.org/>.
- Milner, Danny A. (Jan. 2018). “Malaria Pathogenesis”. en. In: *Cold Spring Harbor Perspectives in Medicine* 8.1, a025569. ISSN: 2157-1422. DOI: 10.1101/cshperspect.a025569. URL: <http://perspectivesinmedicine.cshlp.org/lookup/doi/10.1101/cshperspect.a025569> (visited on 03/24/2023).
- Naslidnyk, Masha et al. (2024). *Comparing Scale Parameter Estimators for Gaussian Process Interpolation with the Brownian Motion Prior: Leave-One-Out Cross Validation and Maximum Likelihood*. arXiv: 2307.07466 [math.ST].
- Price, R.N. et al. (2020). “Plasmodium vivax in the Era of the Shrinking P. falciparum Map”. English. In: *Trends in Parasitology* 36.6, pp. 560–570. ISSN: 1471-4922. DOI: 10.1016/j.pt.2020.03.009.
- Rasmussen, Carl Edward and Christopher K. I. Williams (2008). *Gaussian processes for machine learning*. 3. print. Adaptive computation and machine learning. Cambridge, Mass.: MIT Press. 248 pp. ISBN: 978-0-262-18253-9.
- Smith, David L. et al. (Apr. 2012). “Ross, Macdonald, and a Theory for the Dynamics and Control of Mosquito-Transmitted Pathogens”. en. In: *PLOS Pathogens* 8.4. Publisher: Public Library of Science, e1002588. ISSN: 1553-7374. DOI: 10.1371/journal.ppat.1002588. URL: <https://journals.plos.org/plospathogens/article?id=10.1371/journal.ppat.1002588> (visited on 03/28/2023).
- Srinivas, Niranjana et al. (2010). “Gaussian process optimization in the bandit setting: no regret and experimental design”. In: *Proceedings of the 27th International Conference on International Conference on Machine Learning*. ICML’10. Haifa, Israel: Omnipress, pp. 1015–1022. ISBN: 9781605589077.
- White, Michael T. et al. (Mar. 30, 2016). “Variation in relapse frequency and the transmission potential of *Plasmodium vivax* malaria”. In: *Proceedings of the Royal Society B: Biological Sciences* 283.1827, p. 20160048. ISSN: 0962-8452, 1471-2954. DOI: 10.1098/rspb.2016.0048. URL: <https://royalsocietypublishing.org/doi/10.1098/rspb.2016.0048> (visited on 08/22/2023).
- Wikipedia contributors (2024). *Exponential family* — *Wikipedia, The Free Encyclopedia*. [Online; accessed 16-May-2024]. URL: https://en.wikipedia.org/w/index.php?title=Exponential_family&oldid=1202463189.
- World Health Organization (Dec. 2022). *World malaria report 2022*. en. Tech. rep. Geneva: World Health Organization.
- Zekar, Lara and Tariq Sharman (2023). “Plasmodium Falciparum Malaria”. eng. In: *StatPearls*. Treasure Island (FL): StatPearls Publishing. URL: <http://www.ncbi.nlm.nih.gov/books/NBK555962/> (visited on 03/24/2023).
- Zha, Wen-ting et al. (2020). “Research about the optimal strategies for prevention and control of varicella outbreak in a school in a central city of China: based on an SEIR dynamic model”. en. In: *Epidemiology and Infection* 148, e56. ISSN: 0950-2688, 1469-4409. DOI: 10.1017/S0950268819002188. URL: https://www.cambridge.org/core/product/identifier/S0950268819002188/type/journal_article (visited on 03/26/2024).