

Bayesian Optimisation for Likelihood Free Inference

Make model parameterisation go brrr

Jacob Cumming

University of Melbourne

April 2024



Notation

- ▶ Model is considered a (random) function $f(\boldsymbol{\theta})$ that maps $\boldsymbol{\theta}$ (a vector of parameters) to a model output, that can be transformed into \mathbf{X} , that has the same shape as:
- ▶ \mathbf{X}_{obs} , a vector of outputs given to us usually in the forms of summary statistics (incidence, prevalence, hospitalisations etc).

Parameter inference would become easy if we had...

- ▶ An explicit form for the likelihood: $\mathcal{L}(\theta|\mathbf{X}_{\text{obs}}) := \Pr(\mathbf{X}_{\text{obs}}|\theta)$

Parameter inference would become easy if we had...

- ▶ An explicit form for the likelihood: $\mathcal{L}(\theta|\mathbf{X}_{\text{obs}}) := \Pr(\mathbf{X}_{\text{obs}}|\theta)$
- ▶ Or even $\mathcal{L}(\theta|S(\mathbf{X}_{\text{obs}})) := \Pr(S(\mathbf{X}_{\text{obs}})|\theta)$, where $S(\mathbf{X}_{\text{obs}})$ is a (vector of) summary statistic(s)

Parameter inference would become easy if we had...

- ▶ An explicit form for the likelihood: $\mathcal{L}(\theta|\mathbf{X}_{\text{obs}}) := \Pr(\mathbf{X}_{\text{obs}}|\theta)$
- ▶ Or even $\mathcal{L}(\theta|S(\mathbf{X}_{\text{obs}})) := \Pr(S(\mathbf{X}_{\text{obs}})|\theta)$, where $S(\mathbf{X}_{\text{obs}})$ is a (vector of) summary statistic(s)
- ▶ $\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S(\mathbf{X}_{\text{obs}}))$

Parameter inference would become easy if we had...

- ▶ An explicit form for the likelihood: $\mathcal{L}(\theta|\mathbf{X}_{\text{obs}}) := \Pr(\mathbf{X}_{\text{obs}}|\theta)$
- ▶ Or even $\mathcal{L}(\theta|S(\mathbf{X}_{\text{obs}})) := \Pr(S(\mathbf{X}_{\text{obs}})|\theta)$, where $S(\mathbf{X}_{\text{obs}})$ is a (vector of) summary statistic(s)
- ▶ $\hat{\theta} = \arg \max_{\theta} \mathcal{L}(\theta|S(\mathbf{X}_{\text{obs}}))$
- ▶ $\Pr(\theta|S(\mathbf{X}_{\text{obs}})) \propto \Pr(S(\mathbf{X}_{\text{obs}})|\theta) \Pr(\theta)$

The Sad Truth

- ▶ As models become more complicated, explicit likelihoods don't exist (think agent based models).

A Standard Bayesian Solution

- ▶ Approximate Bayesian Computation (ABC)
 1. Sample from prior
 2. Run model
 3. Accept or reject parameters run based on how well \mathbf{X} 'matches' \mathbf{X}_{obs} .

What is 'matches'

- ▶ Discrepancy function $D : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

- ▶ Can be a norm such as

$$\|S(\mathbf{X}) - S(\mathbf{X}_{\text{obs}})\|_p := (\sum_{i=1}^d |S(\mathbf{X}) - S(\mathbf{X}_{\text{obs}})|^p)^{1/p}$$

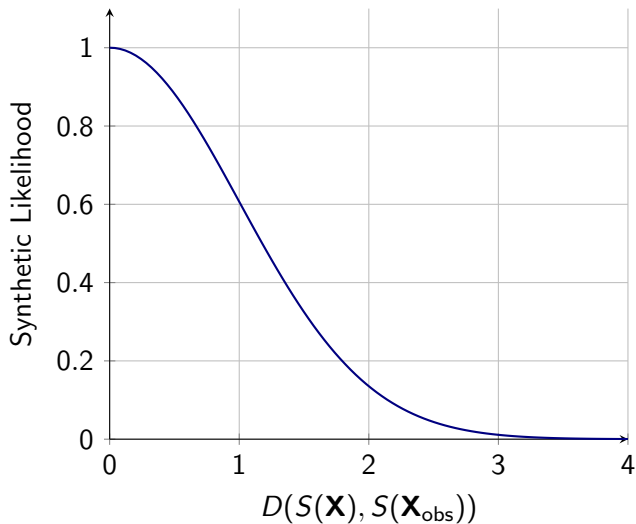
What is 'matches'

- ▶ Discrepancy function $D : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
 - ▶ Can be a norm such as
$$\|S(\mathbf{X}) - S(\mathbf{X}_{\text{obs}})\|_p := (\sum_{i=1}^d |S(\mathbf{X}) - S(\mathbf{X}_{\text{obs}})|^p)^{1/p}$$
 - ▶ Care should be taken to rescale $S(\mathbf{X}_{\text{obs}})$ and $S(\mathbf{X})$ appropriately (ie via a covariance matrix).

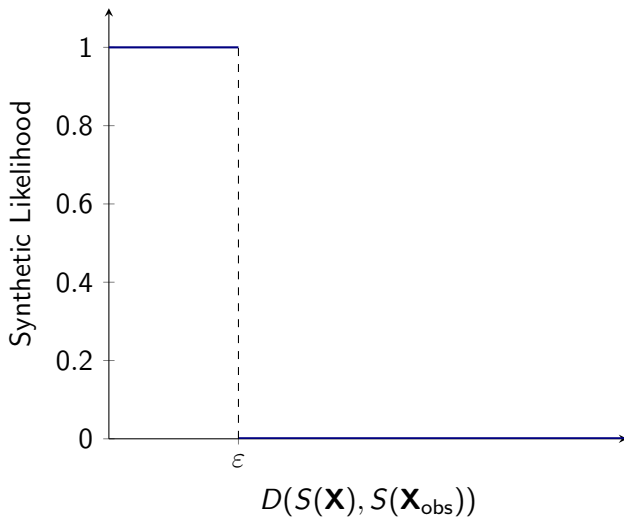
What is 'matches'

- ▶ Discrepancy function $D : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
 - ▶ Can be a norm such as
$$\|S(\mathbf{X}) - S(\mathbf{X}_{\text{obs}})\|_p := (\sum_{i=1}^d |S(\mathbf{X}) - S(\mathbf{X}_{\text{obs}})|^p)^{1/p}$$
 - ▶ Care should be taken to rescale $S(\mathbf{X}_{\text{obs}})$ and $S(\mathbf{X})$ appropriately (ie via a covariance matrix).
- ▶ $D(S(\mathbf{X}), S(\mathbf{X}_{\text{obs}}))$, gives acceptance probability of θ .

Acceptance Probability



Attempt 2



Overall Idea of my Research

- ▶ What if we could 'predict' discrepancy values we hadn't seen before?

Overall Idea of my Research

- ▶ What if we could 'predict' discrepancy values we hadn't seen before?
- ▶ For parameters 'close' to parameters we've already tried it should be easy.

Gaussian Processes

- ▶ Random functions
- ▶ Common examples - Brownian motion, Ornstein Uhlenbeck process

Gaussian Processes on \mathbb{R}^d

Definition (Gaussian Process)

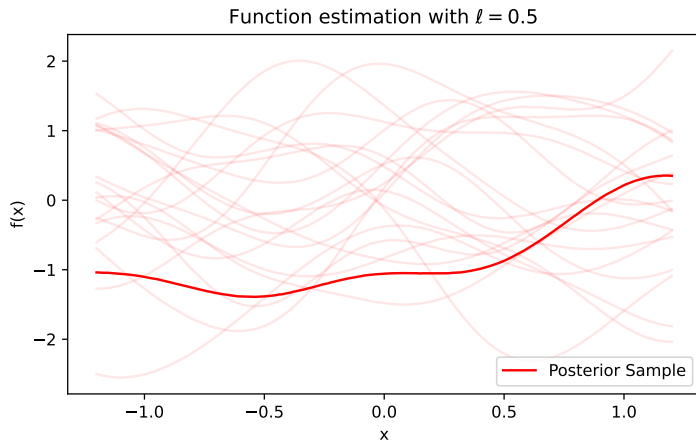
A collection of random variables $\{f(\mathbf{x})\}_{\mathbf{x} \in \mathbb{R}^d}$ is a Gaussian process if all finite dimensional distributions are multivariate normal distributed. That is, there is a function $m : \mathcal{X} \rightarrow \mathbb{R}$ and kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that for all finite sets $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$,

$$\begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix} \sim \text{MVN} \left(\begin{bmatrix} m(\mathbf{x}_1) \\ m(\mathbf{x}_2) \\ \vdots \\ m(\mathbf{x}_n) \end{bmatrix}, \mathbf{K} \right)$$

where

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ k(\mathbf{x}_2, \mathbf{x}_1) & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Gaussian Process Example Realisations



Covariance Kernel Motivation

- ▶ Kernel determines the amount of covariance between sets of indices.
- ▶ When the distance between indices is small, covariance needs to be large

Common Covariance Kernels

- ▶ Matern Kernel

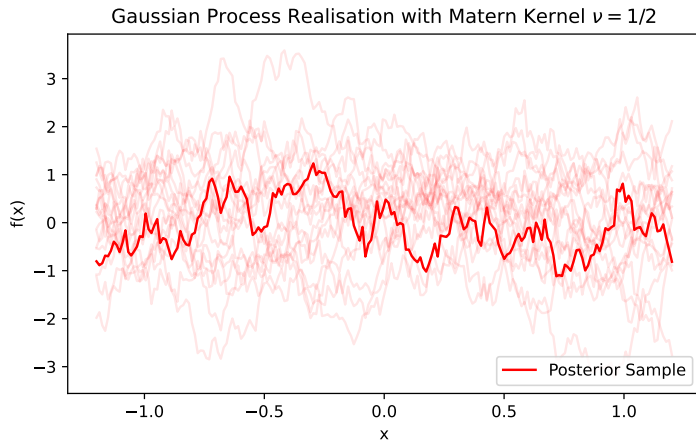
$$k_{\nu}(x, x') = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|x - x'\|}{\ell} \right)^{\nu} K_{\nu} \left(-\frac{\sqrt{2\nu} \|x - x'\|}{\ell} \right)$$

where K_{ν} is a modified Bessel function ($\|\cdot\|$ is the euclidean distance)

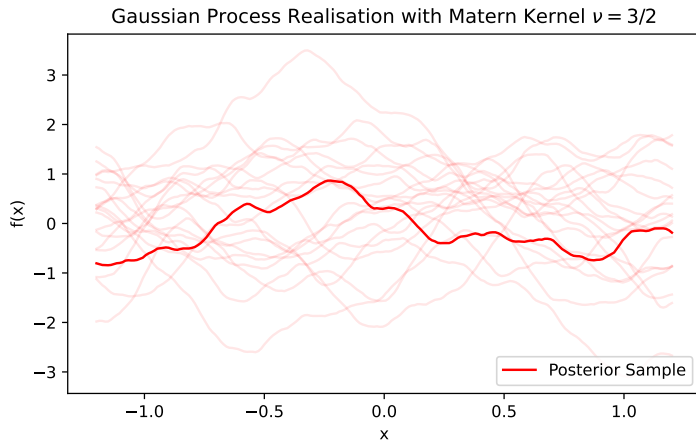
- ▶ $\lfloor \nu \rfloor$ times mean square differentiable.
- ▶ As $\nu \rightarrow \infty$ you get squared exponential covariance function, which results in realisations that are infinitely mean square differentiable:

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

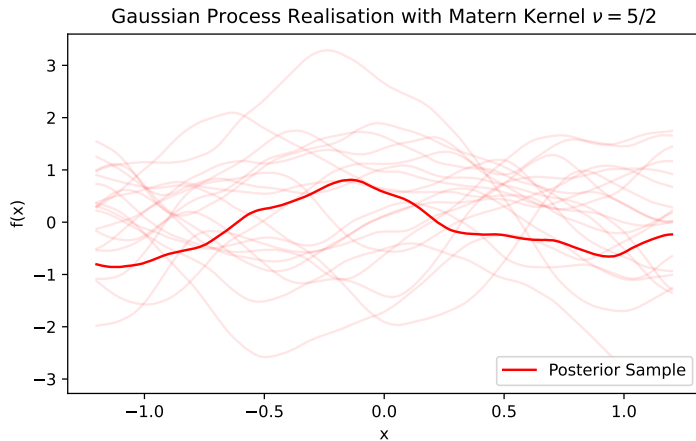
Kernel Choices - Kernel Type



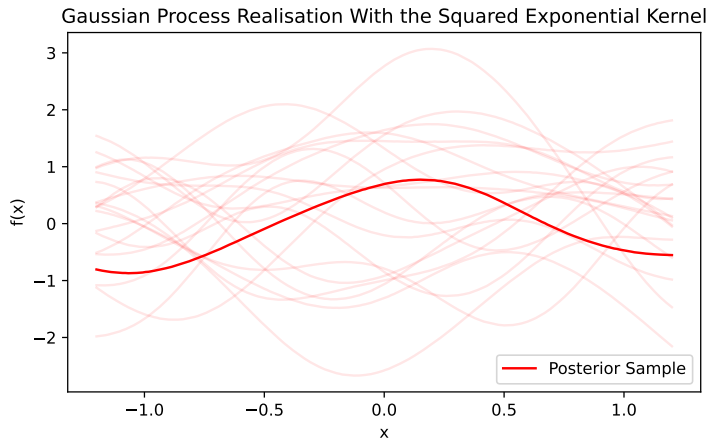
Kernel Choices - Kernel Type



Kernel Choices - Kernel Type



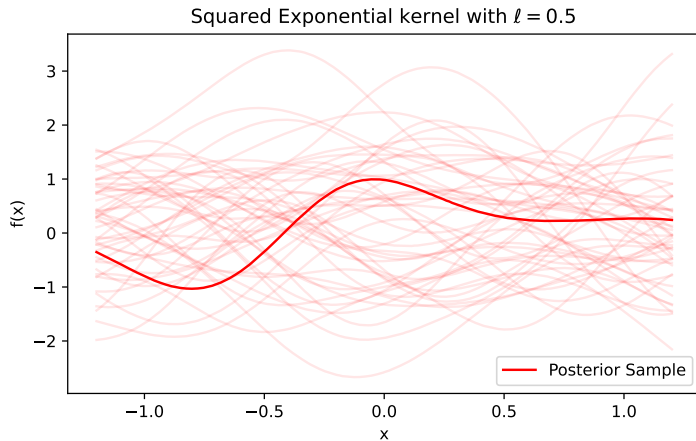
Kernel Choices - Kernel Type



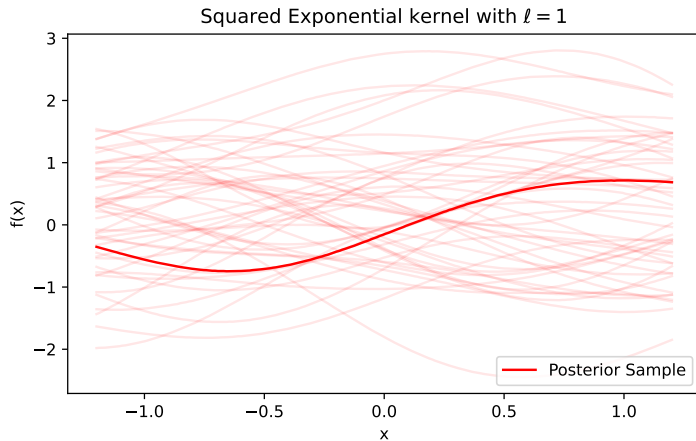
Kernel Choices - length scale

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

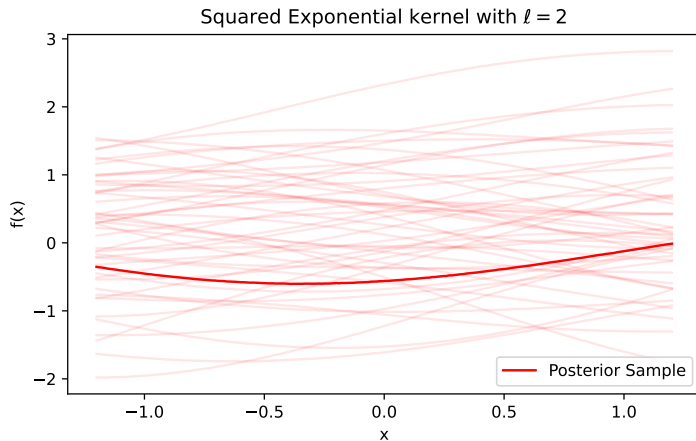
Kernel Choices - length scale



Kernel Choices - length scale



Kernel Choices - length scale

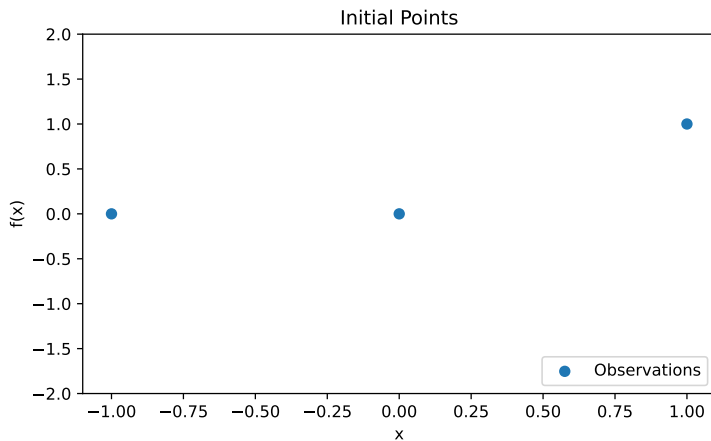


Fitting our GP to data

GPs are 'priors'

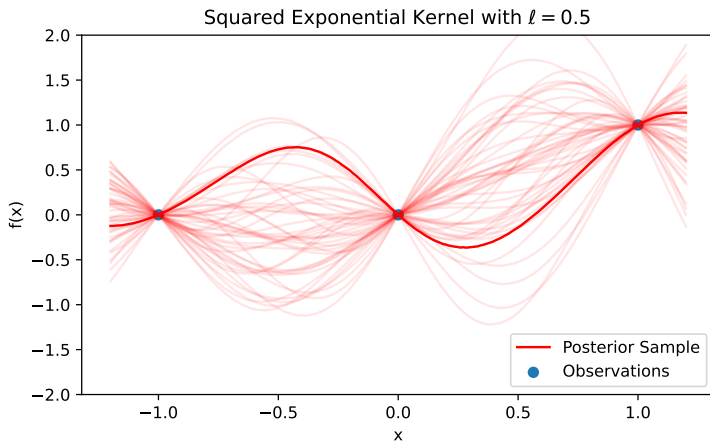
Fitting our GP to data

GPs are 'priors'



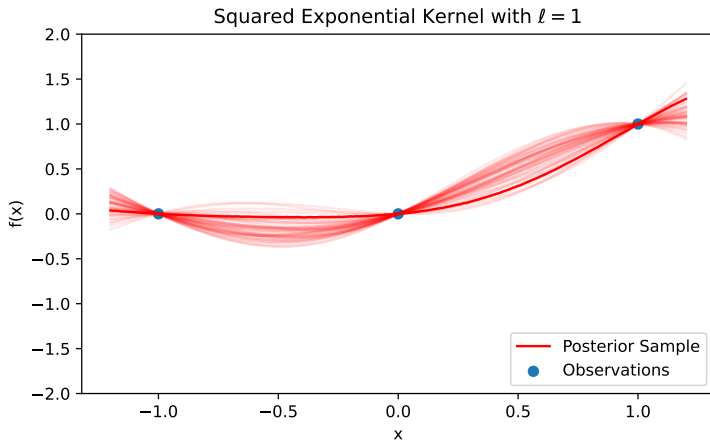
Fitting our GP to data

GPs are 'priors'



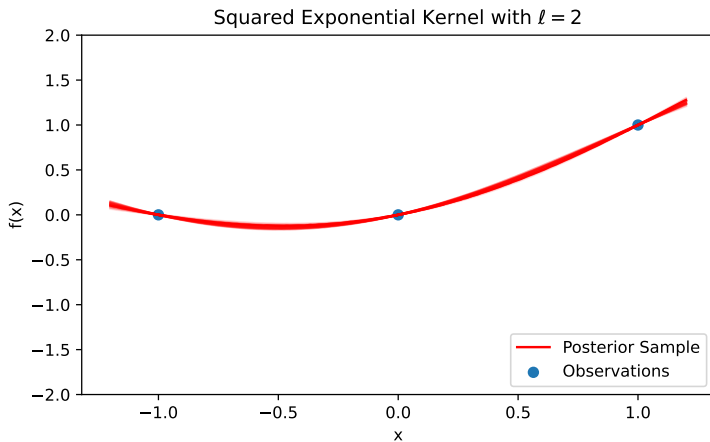
Fitting our GP to data

GPs are 'priors'



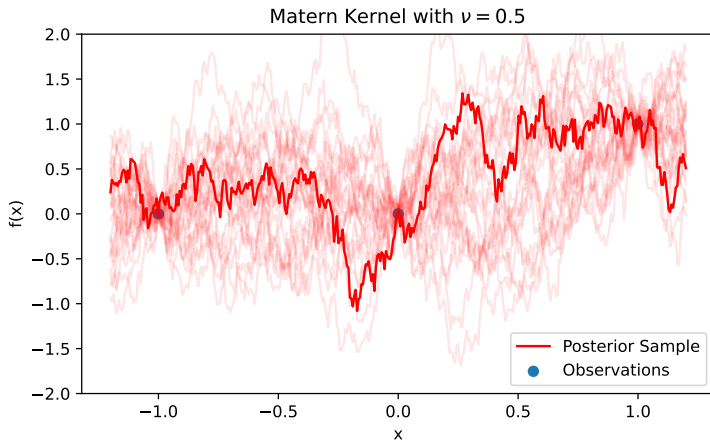
Fitting our GP to data

GPs are 'priors'



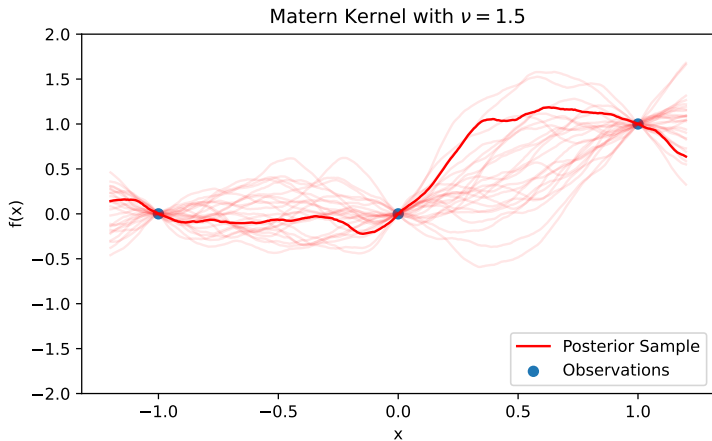
Fitting our GP to data

GPs are 'priors'



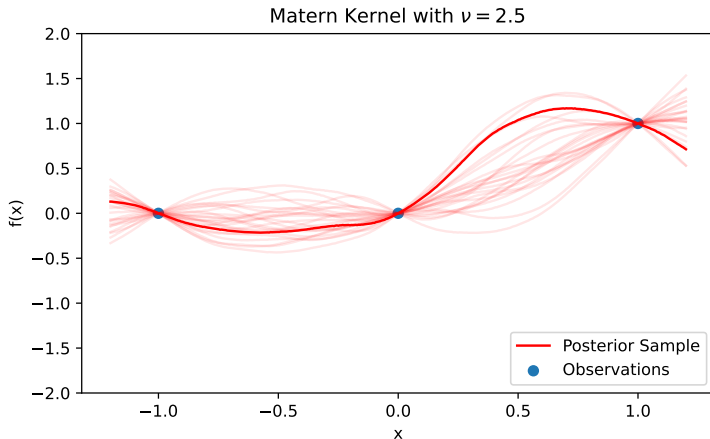
Fitting our GP to data

GPs are 'priors'



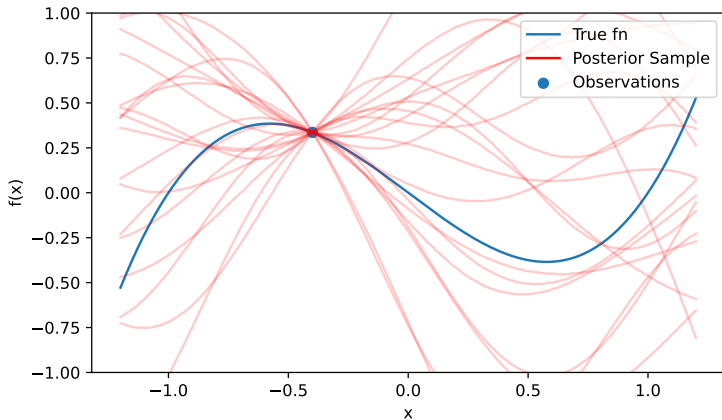
Fitting our GP to data

GPs are 'priors'

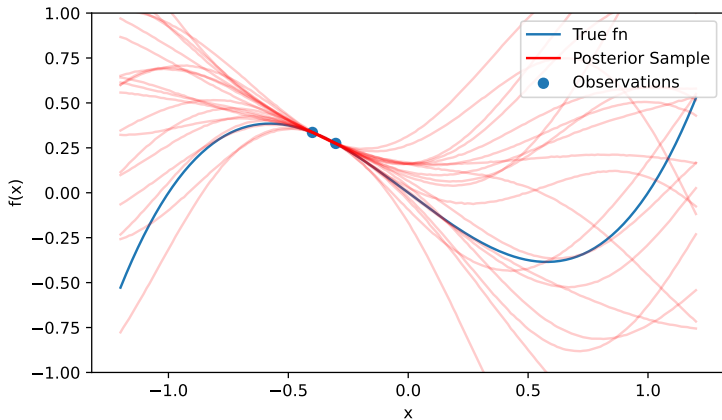


GP regression on $x(x-1)(x+1)$

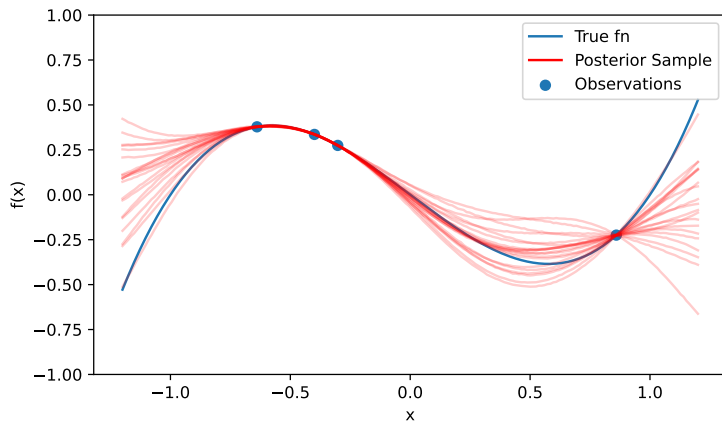
GP regression on $x(x-1)(x+1)$



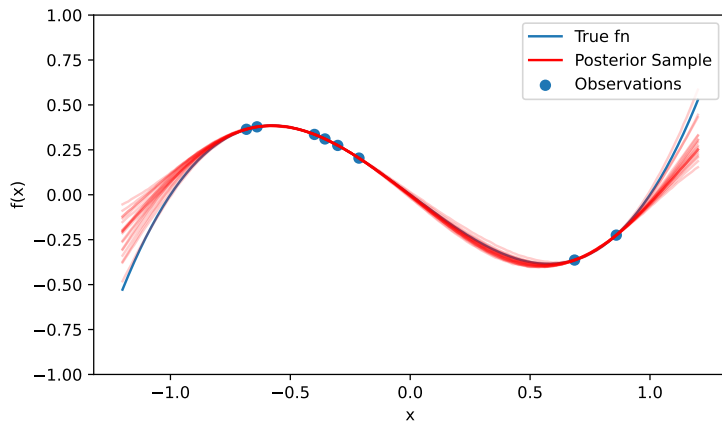
GP regression on $x(x - 1)(x + 1)$



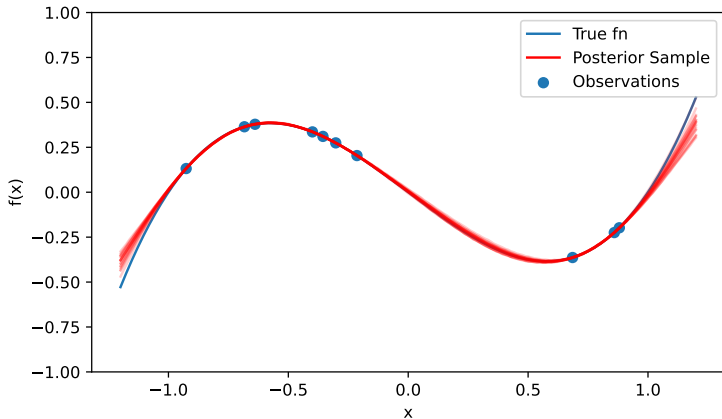
GP regression on $x(x-1)(x+1)$



GP regression on $x(x-1)(x+1)$



GP regression on $x(x-1)(x+1)$



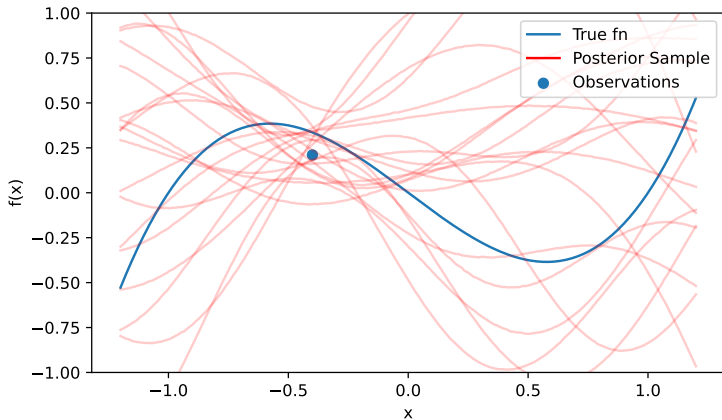
What if we have noise?

Add in observation variance, so that

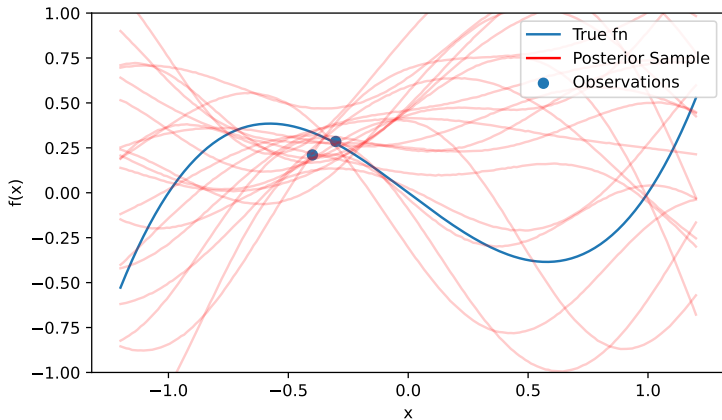
$$\begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix} \sim \text{MVN} \left(\begin{bmatrix} m(\mathbf{x}_1) \\ m(\mathbf{x}_2) \\ \vdots \\ m(\mathbf{x}_n) \end{bmatrix}, \mathbf{K} + \sigma^2 \mathbf{I}_n \right)$$

GP regression on $x(x-1)(x+1)$

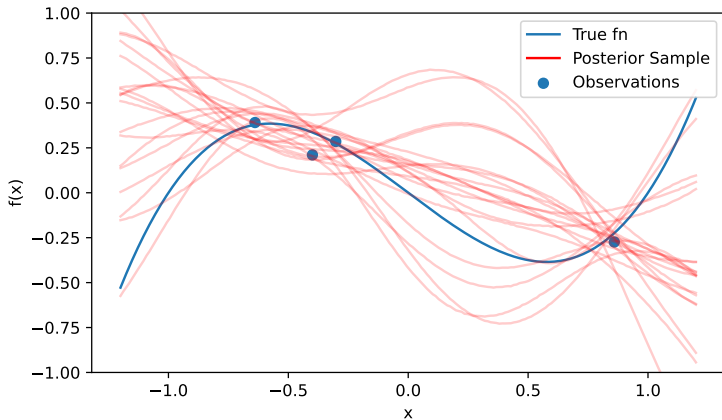
GP regression on $x(x-1)(x+1)$



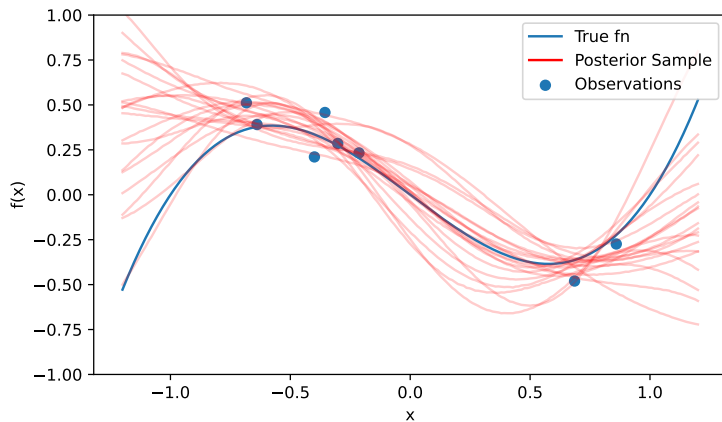
GP regression on $x(x - 1)(x + 1)$



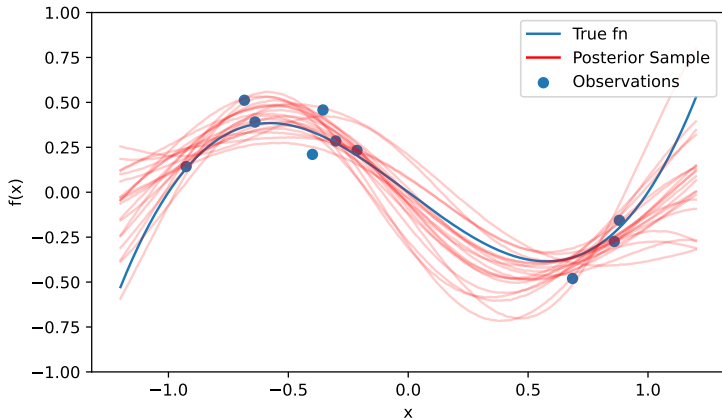
GP regression on $x(x-1)(x+1)$



GP regression on $x(x-1)(x+1)$



GP regression on $x(x-1)(x+1)$



Where to sample function next?

- ▶ Next point $\arg \min_{\theta} A(\theta)$ where A is an acquisition function.

Where to sample function next?

- ▶ Next point $\arg \min_{\theta} A(\theta)$ where A is an acquisition function.
- ▶ BOLFI paper uses

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

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

Where to sample function next?

- ▶ Next point $\arg \min_{\theta} A(\theta)$ where A is an acquisition function.
- ▶ BOLFI paper uses

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

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

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

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

Overall Idea again

- ▶ What if we could 'predict' discrepancy values we hadn't seen before?

Overall Idea again

- ▶ What if we could 'predict' discrepancy values we hadn't seen before?
- ▶ For parameters 'close' to parameters we've already tried it should be easy

Overall Idea again

- ▶ What if we could 'predict' discrepancy values we hadn't seen before?
- ▶ Use Gaussian process to predict discrepancy function

About Vivax Malaria

- ▶ Has dormant liver stage on top of blood stage infection

Champagne Model Parameters

- ▶ α : proportion of those infected but cleared of blood stage infections (through treatment)
- ▶ β : a further proportion that are also cleared of liver stage parasites, given that they were also cleared of blood stage infection (radical cure)
- ▶ λ : the rate of infection
- ▶ γ_L : rate of clearance of liver stage disease
- ▶ f : rate of relapse
- ▶ r : rate of blood stage clearance
- ▶ δ : importation rate (which we assume is 0)

Champagne ODEs

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

$$\frac{dI_0}{dt} = -(\lambda I_{\text{total}} + \delta)I_0 + \gamma_L I_L - rI_0$$

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

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

Champagne Model Diagram

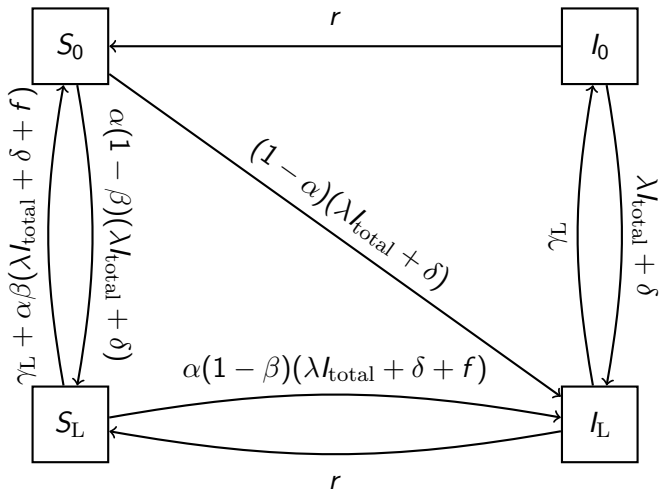


Figure: *P. vivax* model described by Champagne et al. 2022

Model Calibration Data

- ▶ Doob-Gillespie algorithm with paper parameters reported in the original paper for 'observed data', 10 initial infections, 1000 people.

Model Calibration Data

- ▶ Doob-Gillespie algorithm with paper parameters reported in the original paper for 'observed data', 10 initial infections, 1000 people.
- ▶ $S(\mathbf{X}_{\text{obs}}) := \{w_{\text{obs}}, p_{\text{obs}}, m_{\text{obs}}\}$
 - ▶ w_{obs} : weekly incidence around (stochastic) equilibrium
 - ▶ p_{obs} : prevalence around (stochastic) equilibrium
 - ▶ m_{obs} : incidence in the first month of the epidemic

Model Calibration Data

- ▶ Doob-Gillespie algorithm with paper parameters reported in the original paper for 'observed data', 10 initial infections, 1000 people.
- ▶ $S(\mathbf{X}_{\text{obs}}) := \{w_{\text{obs}}, p_{\text{obs}}, m_{\text{obs}}\}$
 - ▶ w_{obs} : weekly incidence around (stochastic) equilibrium
 - ▶ p_{obs} : prevalence around (stochastic) equilibrium
 - ▶ m_{obs} : incidence in the first month of the epidemic
- ▶ $D(S(\mathbf{X}), S(\mathbf{X}_{\text{obs}})) := \sqrt{\left(\frac{w_{\text{obs}} - w}{w_{\text{obs}}}\right)^2 + \left(\frac{p_{\text{obs}} - p}{p_{\text{obs}}}\right)^2 + \left(\frac{p_{\text{obs}} - p}{p_{\text{obs}}}\right)^2}$
 - ▶ L_2 norm on the relative differences

Big Problems (Big Solutions?)

- ▶ Observation variance is considered constant across the GP

Big Problems (Big Solutions?)

- ▶ Observation variance is considered constant across the GP
 - ▶ Fix by modelling observation variance as another GP

Big Problems (Big Solutions?)

- ▶ Observation variance is considered constant across the GP
 - ▶ Fix by modelling observation variance as another GP
- ▶ Observation variance sometimes falls off a cliff