

# Machine Learning

## Lecture 8 - Gaussian process regression and classification

**Mattias Villani**

Department of Statistics  
Stockholm University

Department of Computer and Information Science  
Linköping University



[mattiasvillani.com](http://mattiasvillani.com)



@matvil



mattiasvillani

# Lecture overview

- Bayesian inference
- Gaussian process regression
- Gaussian process classification

# Bayesian inference

- Parametric model  $p(x|\theta)$ .
- **Likelihood**:  $p(x_1, \dots, x_n|\theta)$ .
- Bayesian inference uses Bayes' theorem to combine
  - ▶ data information (likelihood)
  - ▶ other information (prior)
- **Prior** distribution  $p(\theta)$ .
- **Subjective probability** for **unknown** quantities.
- **Posterior distribution**

$$p(\theta|x_1, \dots, x_n) \propto p(x_1, \dots, x_n|\theta)p(\theta)$$

$$\text{Posterior} \propto \text{Likelihood} \times \text{Prior}$$

# Normal data, known variance - normal prior

## ■ Model

$$x_1, \dots, x_n | \theta, \sigma^2 \stackrel{\text{iid}}{\sim} N(\theta, \sigma^2).$$

## ■ Prior

$$\theta \sim N(\mu_0, \tau_0^2)$$

## ■ Posterior

$$\begin{aligned} p(\theta | x_1, \dots, x_n) &\propto p(x_1, \dots, x_n | \theta, \sigma^2) p(\theta) \\ &\propto N(\theta | \mu_n, \tau_n^2), \end{aligned}$$

where

$$\frac{1}{\tau_n^2} = \frac{n}{\sigma^2} + \frac{1}{\tau_0^2},$$

$$\mu_n = w\bar{x} + (1 - w)\mu_0,$$

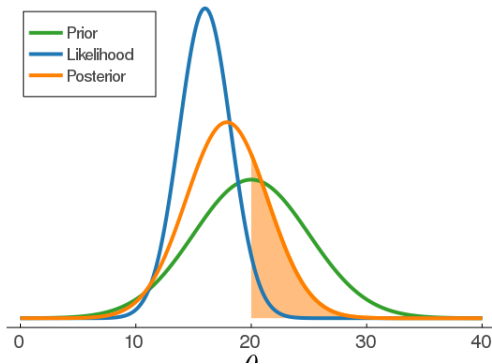
and

$$w = \frac{\frac{n}{\sigma^2}}{\frac{n}{\sigma^2} + \frac{1}{\tau_0^2}}.$$

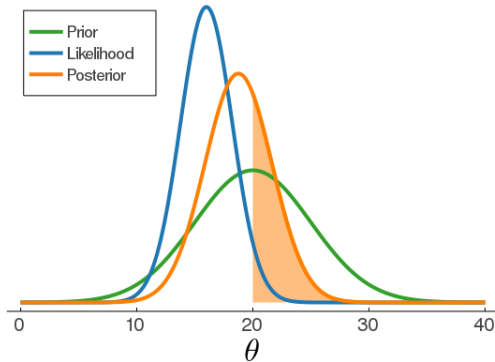
# Download speed

- **Problem:** My internet provider promises an average download speed of at least 20 Mbit/sec. Are they lying?
- **Data:**  $x = (22.42, 34.01, 35.04, 38.74, 25.15)$  Mbit/sec.
- **Model:**  $X_1, \dots, X_5 \sim N(\theta, \sigma^2)$ .
- Assume  $\sigma = 5$  (measurements can vary  $\pm 10$  MBit with 95% probability)
- My **prior:**  $\theta \sim N(20, 5^2)$ .

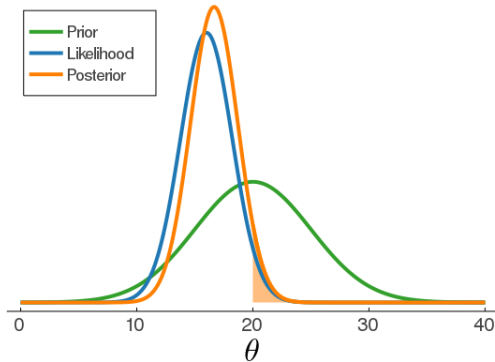
# Download speed $n=1$



# Download speed $n=2$

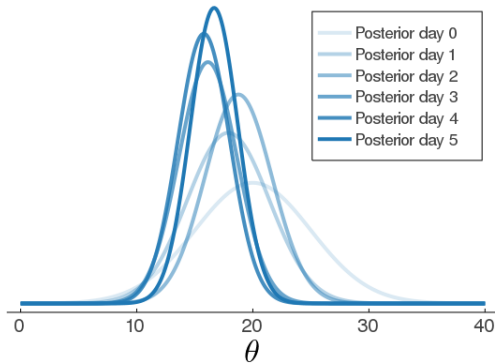


# Download speed $n=5$





# Bayesian updating



# Linear regression with known variance

- The linear regression model in **matrix form**

$$\underset{(n \times 1)}{y} = \underset{(n \times k)}{X} \underset{(k \times 1)}{\beta} + \underset{(n \times 1)}{\varepsilon}, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

- **Prior** for  $\beta$

$$\beta \sim N(0, \sigma^2 \Omega_0^{-1})$$

- **Posterior**

$$\beta | y, X \sim N[\mu_n, \sigma^2 \Omega_n^{-1}]$$

$$\mu_n = (X^\top X + \Omega_0)^{-1} X^\top X \hat{\beta}$$

$$\Omega_n = X^\top X + \Omega_0$$

- **Posterior mean** estimate  $\mu_n$  is a shrunken version of least squares estimate  $\hat{\beta} = (X^\top X)^{-1} X^\top y$ .
- Prior acts as **regularization**.  $\Omega_0 = \lambda I$  gives **Ridge**.

# Nonlinear regression

## ■ Linear regression

$$y = f(x) + \epsilon$$

$$f(x) = x^T \beta$$

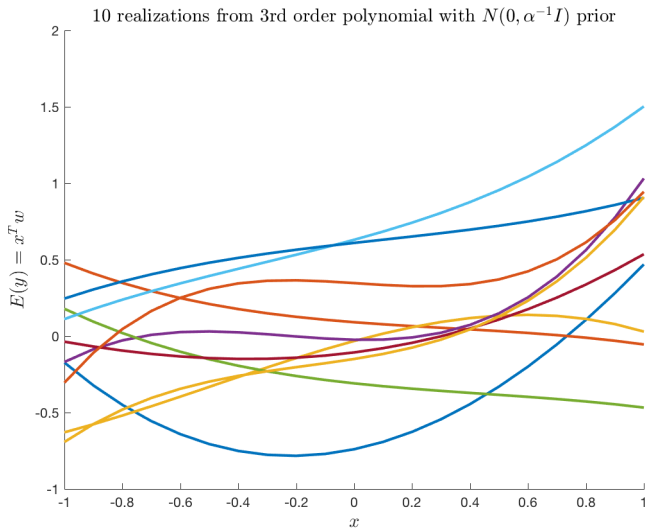
and  $\epsilon \sim N(0, \sigma_n^2)$  and iid over observations.

## ■ Polynomial regression: $\phi(x) = (1, x, x^2, x^3, \dots, x^k)$ :

$$f(x) = \phi(x)^T \beta.$$

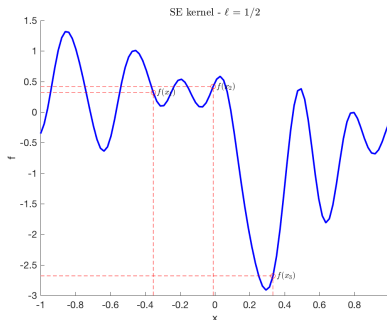
## ■ More generally: **splines** with **basis functions**.

# A prior on $\beta$ is really a prior over functions



# Non-parametric regression

- **Non-parametric regression**: avoid a parametric form for  $f(\cdot)$ .
- Treat  $f(x)$  as **an unknown parameter for every  $x$** .



- A *new* parameter for every  $x$ , you must be joking?
- Instead of restricting to linear, impose **smoothness**.

# Two views on GPs

- **Weight space view**
- Restrict attention to a grid of  $x$ -values:  $x_1, \dots, x_k$ .
- Put a joint prior on the **vector of  $k$  function values**

$$f(x_1), \dots, f(x_k)$$

---

- **Function space view**
- Treat  $f$  as an **unknown function**.
- Put a prior over a set of functions.

# Gaussian process and its kernel

- A GP implies:

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_k) \end{pmatrix} \sim N(m, K)$$

- But how do we specify the  $k \times k$  **covariance matrix**  $K$ ?

$$\text{Cov}(f(x_p), f(x_q))$$

- **Squared exponential covariance function**

$$\text{Cov}(f(x_p), f(x_q)) = k(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2} \left(\frac{x_p - x_q}{\ell}\right)^2\right)$$

- Nearby  $x$ 's have highly correlated function ordinates  $f(x)$ .
- We can compute  $\text{Cov}(f(x_p), f(x_q))$  for *any*  $x_p$  and  $x_q$ .

# Gaussian processes

## Definition

A **Gaussian process (GP)** is a collection of random variables, any finite number of which have a multivariate Gaussian distribution.

- A GP is a **probability distribution over functions**.
- A GP is specified by a **mean** and a **covariance function**

$$m(x) = E[f(x)]$$

$$k(x, x') = E[(f(x) - m(x))(f(x') - m(x'))]$$

for any two inputs  $x$  and  $x'$ .

- A **Gaussian process** is denoted by

$$f(x) \sim \text{GP}(m(x), k(x, x'))$$

- $f(x) \sim \text{GP}$  encodes **prior beliefs** about the unknown  $f(\cdot)$ .



# Gaussian processes

■ Let  $r = \|x - x'\|$ .

■ **Squared exponential (SE)** kernel ( $\ell > 0, \sigma_f > 0$ )

$$K_{SE}(r) = \sigma_f^2 \exp\left(-\frac{r^2}{2\ell^2}\right)$$

■ **Matérn** kernel ( $\ell > 0, \sigma_f > 0, \nu > 0$ )

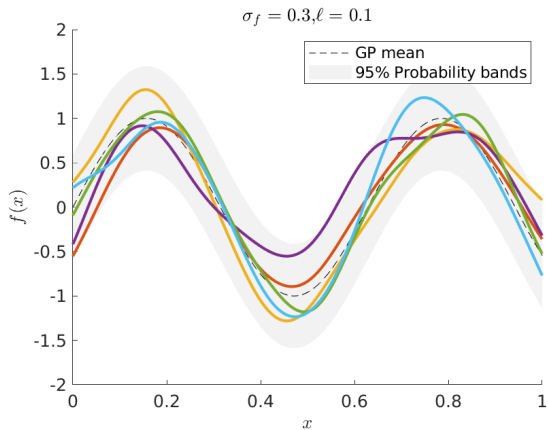
$$K_{Matern}(r) = \sigma_f^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu}r}{\ell}\right)^\nu K_\nu\left(\frac{\sqrt{2\nu}r}{\ell}\right)$$

■ **Simulate draw** from  $f(x) \sim \text{GP}(m(x), k(x, x'))$  by:

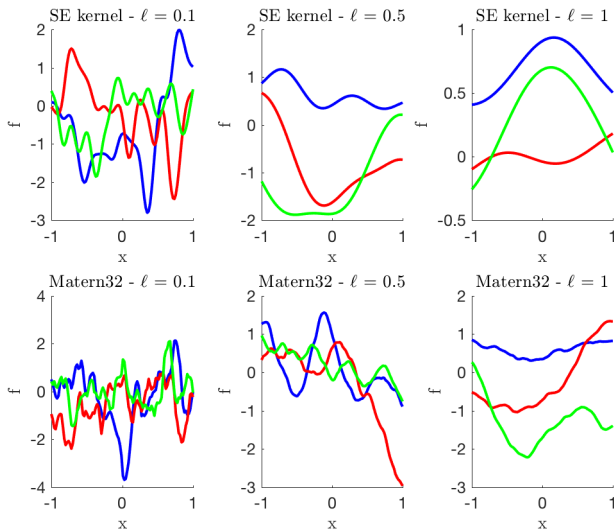
- ▶ form a grid  $x_* = (x_1, \dots, x_n)$
- ▶ simulate function values from multivariate normal:

$$f(x_*) \sim N(m(x_*), K(x_*, x_*))$$

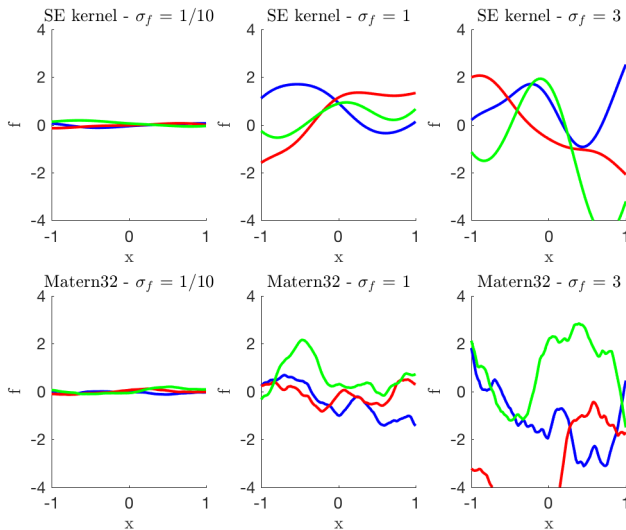
# Simulating a GP



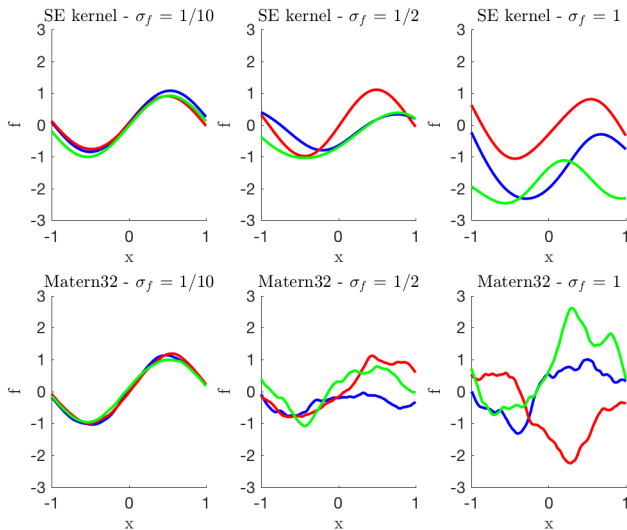
# The length scale $\ell$ determines the smoothness



# The scale factor $\sigma_f$ determines the variance



The mean can be  $\sin(3x)$ . Or whatever.



# Sequential simulation of GPs

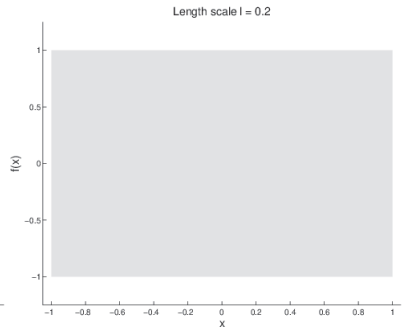
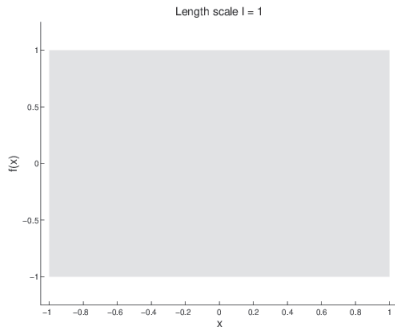
- The joint way: Choose a grid  $x_1, \dots, x_k$ . Simulate the  $k$ -vector

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_k) \end{pmatrix} \sim N(m, K)$$

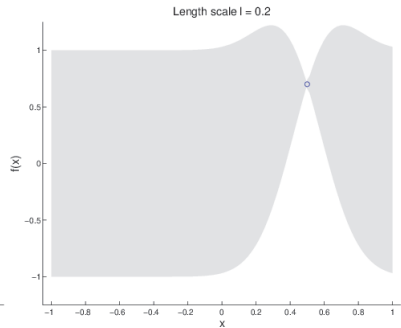
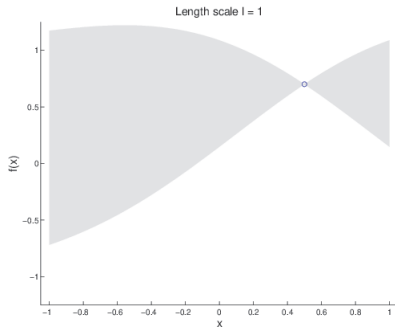
- More intuition from the conditional decomposition

$$\begin{aligned} p(f(x_1), f(x_2), \dots, f(x_k)) &= p(f(x_1)) p(f(x_2)|f(x_1)) \cdots \\ &\quad \times p(f(x_k)|f(x_1), \dots, f(x_{k-1})) \end{aligned}$$

# Simulating from $p(f(x_1))$

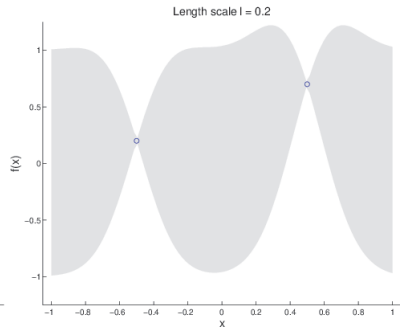
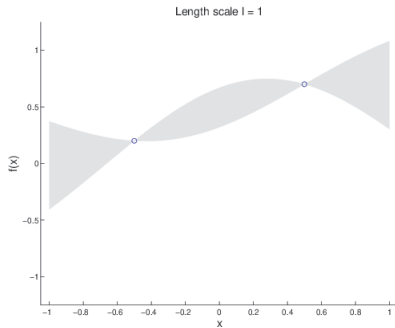


# Simulating from $p(f(x_2)|f(x_1))$

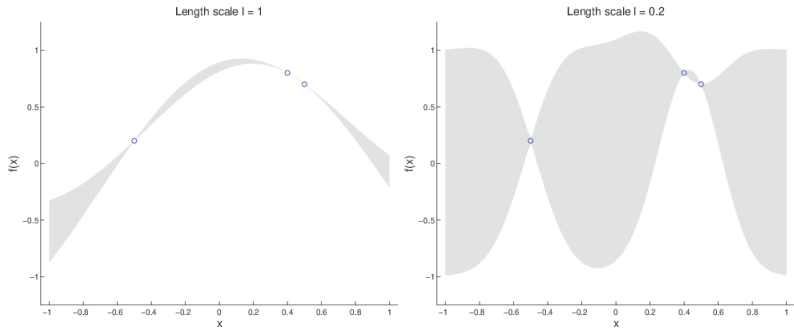




# Simulating from $p(f(x_3)|f(x_1), f(x_2))$



# Simulating from $p(f(x_4)|f(x_1), f(x_2), f(x_3))$



# The posterior for a Gaussian Process Regression

## Model

$$y_i = f(x_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma_n^2)$$

## Prior

$$f(x) \sim GP(0, k(x, x'))$$

Observed:  $x = (x_1, \dots, x_n)^T$  and  $y = (y_1, \dots, y_n)^T$ .

Goal: posterior of  $f(\cdot)$  over a grid of  $x$ -values:  $f_* = f(x_*)$ .

## Posterior

$$f_* | x, y, x_* \sim N(\bar{f}_*, \text{cov}(f_*))$$

$$\bar{f}_* = K(x_*, x) [K(x, x) + \sigma_n^2 I]^{-1} y$$

$$\text{cov}(f_*) = K(x_*, x_*) - K(x_*, x) [K(x, x) + \sigma_n^2 I]^{-1} K(x, x_*)$$

# Sketch for proof of posterior

- Idea: obtain joint  $p(y, f_*)$  and then  $p(f_*|y)$  by conditioning.

- **Model**

$$y_i = f(x_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma_n^2)$$

- **Prior**

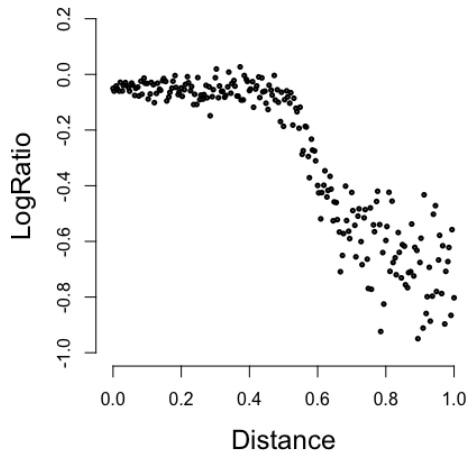
$$f(x) \sim GP(0, k(x, x'))$$

- Joint distribution of  $(y, f_*)$

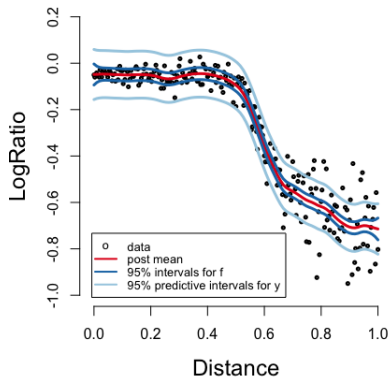
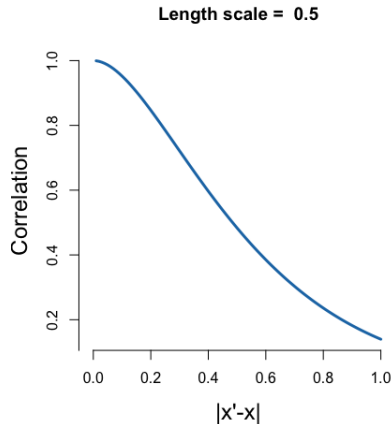
$$\begin{pmatrix} y \\ f_* \end{pmatrix} \sim N \left[ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} K(x, x) + \sigma_n^2 I & K(x, x_*) \\ K(x_*, x) & K(x_*, x_*) \end{pmatrix} \right]$$

- Result: conditional distributions from multivariate normal are normal.

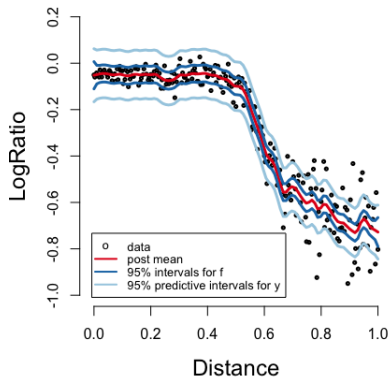
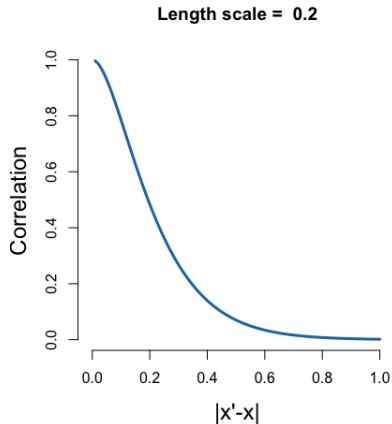
## Example - LIDAR data



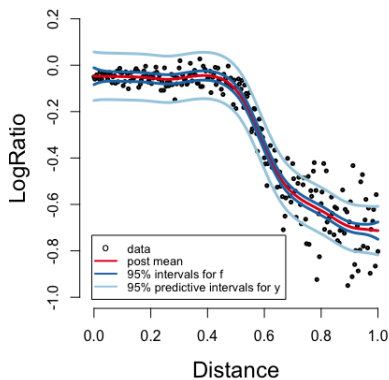
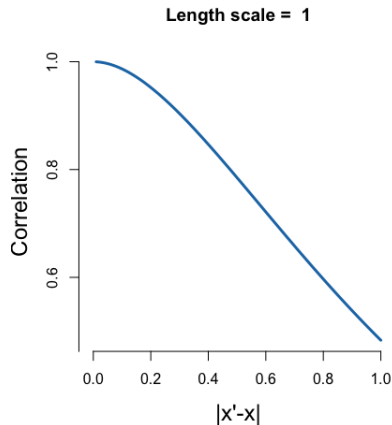
# GP fit to LIDAR data $\ell = 0.5, \sigma_f = 0.5, \sigma_n = 0.05$



# GP fit to LIDAR data $\ell = 0.2, \sigma_f = 0.5, \sigma_n = 0.05$

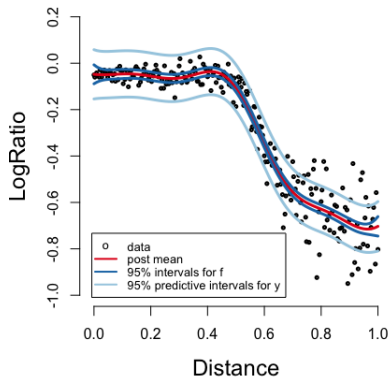
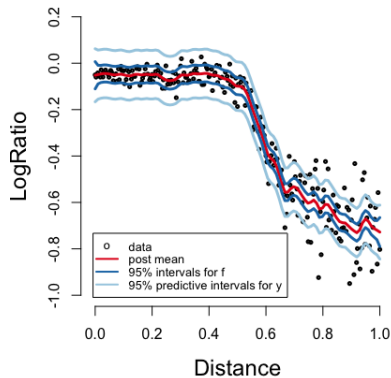


# GP fit to LIDAR data $\ell = 1, \sigma_f = 0.5, \sigma_n = 0.05$





# Matern32 vs SquaredExp for $\ell = 0.2$



# Inference for the hyperparameters

- Kernel depends on **hyperparameters**  $\theta = (\sigma_f, \ell)^T$ . Example

$$k(x, x') = \sigma_f^2 \exp \left( -\frac{1}{2} \frac{\|x - x'\|^2}{\ell^2} \right)$$

- Common: maximize the **marginal likelihood** wrt  $\theta$ :

$$p(y|X, \theta) = \int p(y|X, f, \theta) p(f|X, \theta) df$$

$f = f(X)$  is a vector of function values in the training data.

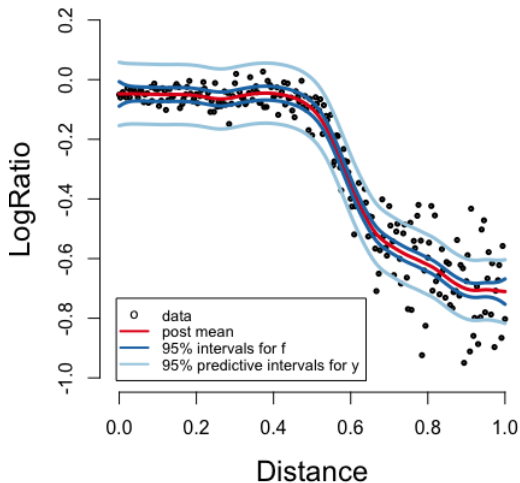
- For **Gaussian process regression**:

$$\log p(y|X, \theta) = -\frac{1}{2} y^T (K + \sigma_n^2 I)^{-1} y - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{n}{2} \log(2\pi)$$

- Proper **Bayesian inference** for hyperparameters

$$p(\theta|y, X) \propto p(y|X, \theta) p(\theta).$$

# GP fit LIDAR $\ell_{opt} = 0.61, \sigma_{f,opt} = 0.44, \sigma_n = 0.05$



# GP Classification

- **Binary** or multi-class **response**. Aim:  $\Pr(y_i = 1|x_i)$ .
- **Logistic regression**

$$\Pr(y_i = 1|x_i) = \lambda(x_i^T \beta), \text{ where } \lambda(z) = \frac{1}{1 + \exp(-z)}.$$

- $\lambda(z)$  'squashes' the linear prediction  $x^T \beta \in \mathbb{R}$  into  $[0, 1]$ .
- **Linear decision boundaries** because of linear predictor  $x^T \beta$ .
- **GP classification**: replace  $x^T \beta$  by  $f(x)$  where

$$f \sim \text{GP}(0, k(x, x'))$$

and squash  $f$  through logistic function

$$\Pr(y = 1|x) = \lambda(f(x))$$

- Nonparametric **flexible decision boundaries**.

# GP Classification on simulated data

