

Advanced Bayesian Learning

Gaussian Process Regression and Classification - Lecture 1

Mattias Villani

**Department of Statistics
Stockholm University**

Department of Computer and Information Science
Linköping University



Course overview

■ Four topics

- ▶ Gaussian Process Regression and Classification
- ▶ Bayesian Nonparametrics
- ▶ Variational Inference
- ▶ Bayesian Regularization

■ Examination

- ▶ Individual Lab/Exercise for each topic
- ▶ Deadline for submission: day before new topic starts.
- ▶ Extra deadline for all four topics: Sept 15, 2024.

Topic overview

- Gaussian Process Regression
- Gaussian Process Classification

Nonlinear regression

■ Linear regression

$$y = f(\mathbf{x}) + \epsilon$$

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$$

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

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

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}.$$

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

■ Example: **thin plate splines** with knots $\kappa_1, \dots, \kappa_N$ in \mathbf{x} -space

$$\phi_k(\mathbf{x}) = \ln(\|\mathbf{x} - \kappa_k\|) \|\mathbf{x} - \kappa_k\|^2$$

Recap: Bayesian linear regression

■ Prior

$$w \sim N(0, \Sigma_p)$$

■ Posterior [X is $D \times n$]

$$w|X, y \sim N(\bar{w}, A^{-1})$$

$$A = \sigma_n^{-2} X X^T + \Sigma_p^{-1}$$

$$\bar{w} = \sigma_n^{-2} A^{-1} X y = \left(X X^T + \sigma_n^2 \Sigma_p^{-1} \right)^{-1} X y$$

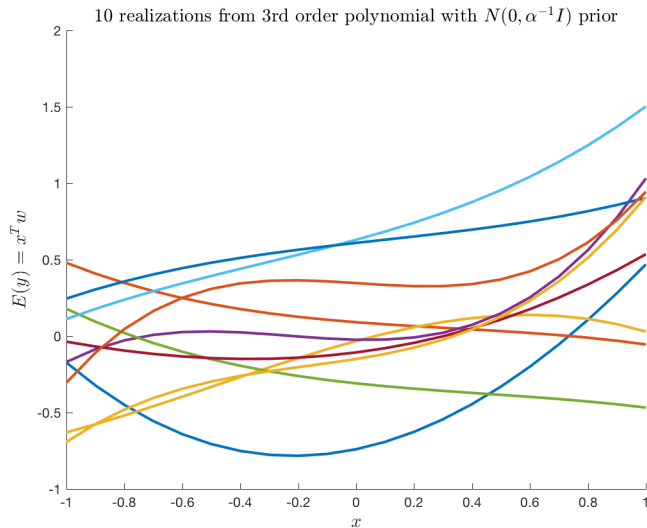
■ Predictive density for mean $f(x_*) = x_*^T w$ at new x_*

$$f(x_*)|x_*, X, y \sim N\left(x_*^T \bar{w}, x_*^T A^{-1} x_*\right)$$

■ Predictive density for new response y_*

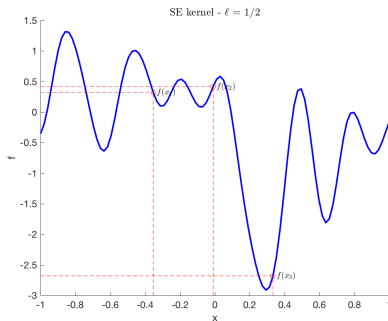
$$y_*|x_*, X, y \sim N\left(x_*^T \bar{w}, x_*^T A^{-1} x_* + \sigma_n^2\right)$$

A prior on w is really a prior over functions



Non-parametric regression

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



- A *new* parameter for every x !
- Instead of restricting to linear, impose “**prior 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 (thank you, Kolmogorov!)

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 ($\ell > 0, \sigma_f > 0$)

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

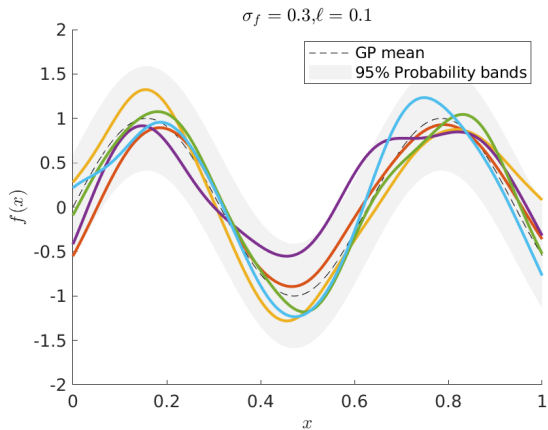
- **Matérn** kernel ($\ell > 0, \sigma_f > 0, \nu > 0$)

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

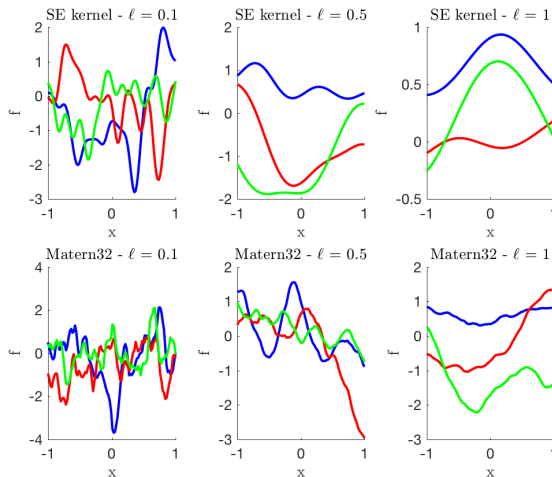
- **Simulate draw** from $f(x) \sim \text{GP}(m(x), k(x, x'))$ by:
 - ▶ form a grid $x_* = (x_1, \dots, x_n)$
 - ▶ simulate function values from multivariate normal:

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

Simulating a GP

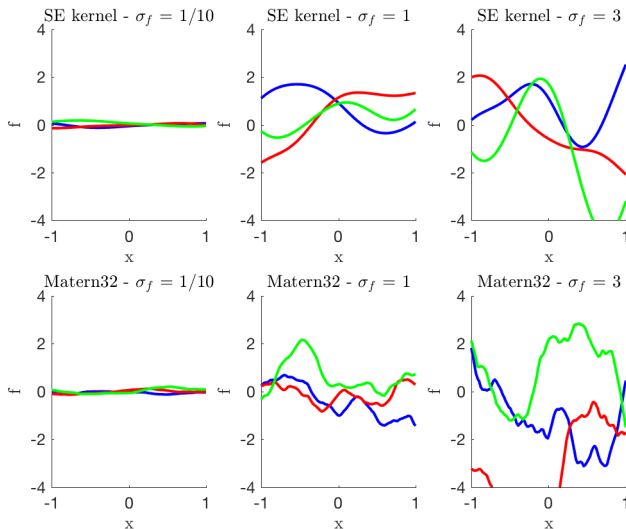


The length scale ℓ - the correlation distance

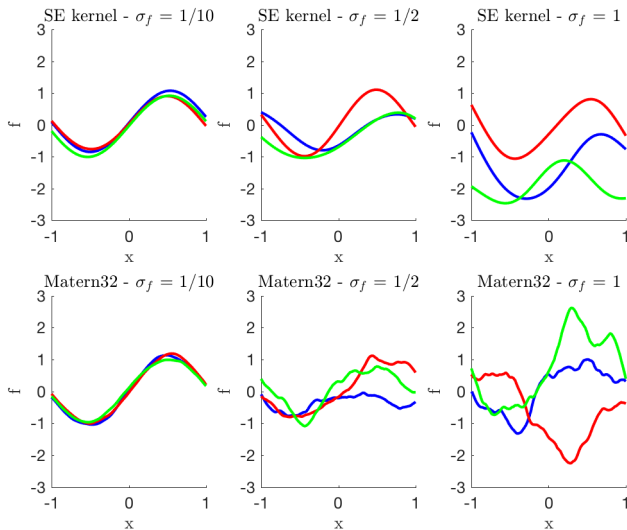


- SE: expected number of zero-crossings on $[0, 1]$: $(2\pi\ell)^{-1}$ (Eq. 4.3)

The scale factor σ_f determines the variance



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



Sequential simulation of GPs

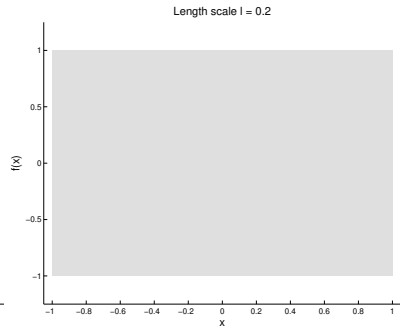
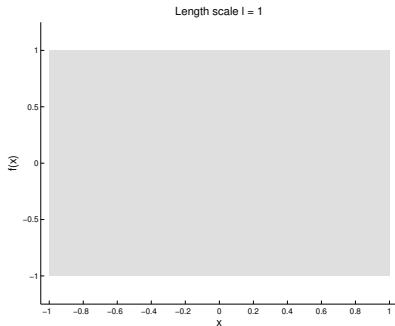
- The joint way: Choose a grid x_1, \dots, x_k . Simulate the k -vector

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

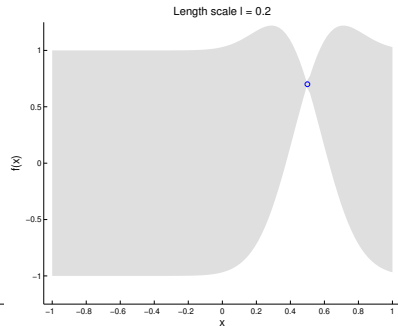
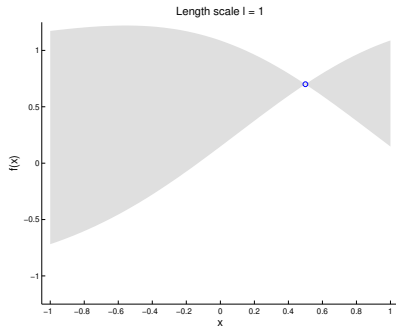
- More intuition from the conditional decomposition

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

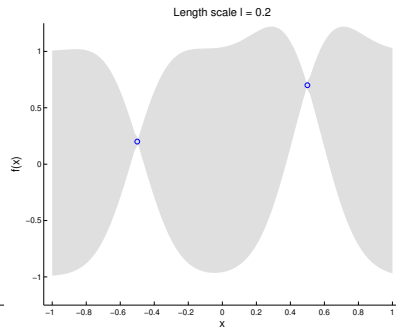
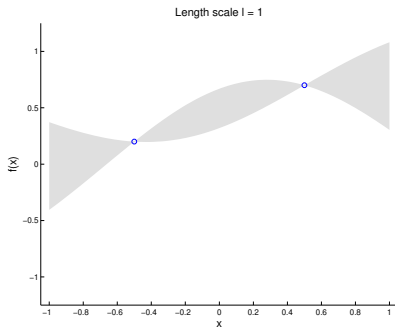
Simulating from $p(f(x_1))$



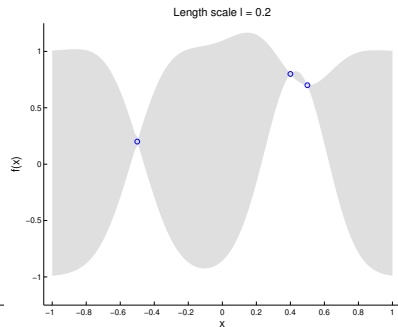
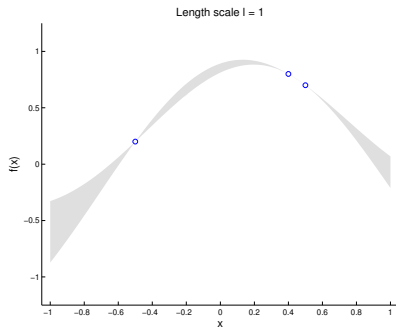
Simulating from $p(f(x_2)|f(x_1))$



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



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



Multivariate normal distribution

- The density of the p -variate $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is

$$f(\mathbf{x}) = \left(\frac{1}{2\pi}\right)^{p/2} \frac{1}{\sqrt{\det \boldsymbol{\Sigma}}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- **Linear combinations.** Let $\mathbf{y} = \mathbf{B}\mathbf{x} + \mathbf{b}$, then

$$\mathbf{y} \sim N(\mathbf{B}\boldsymbol{\mu} + \mathbf{b}, \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}')$$

- Let $\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$ where \mathbf{x}_1 is $p_1 \times 1$ and \mathbf{x}_2 is $p_2 \times 1$ and

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \text{ and } \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}$$

- **Marginals are normal.** Let $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then

$$\mathbf{x}_1 \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

- **Conditionals are normal.** Let $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then

$$\mathbf{x}_2 | \mathbf{x}_1 = \mathbf{x}_1^* \sim N \left[\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x}_1^* - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \right]$$

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 test data: $f_* = f(x_*)$.

■ Posterior

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

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

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

■ **Predictive distribution** for new test data

$$y_* | x, y, x_* \sim N(\bar{f}_*, \text{cov}(f_*) + \sigma_n^2 I)$$

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

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

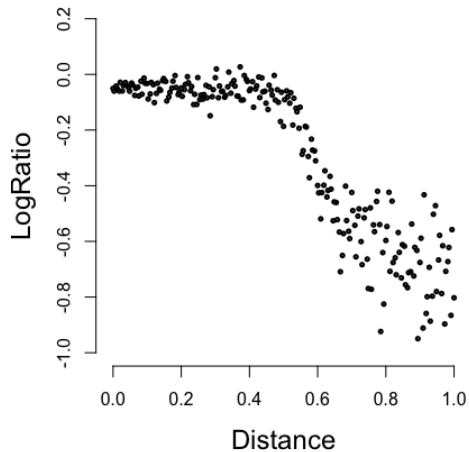
- Joint distribution of (y, f_*)

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

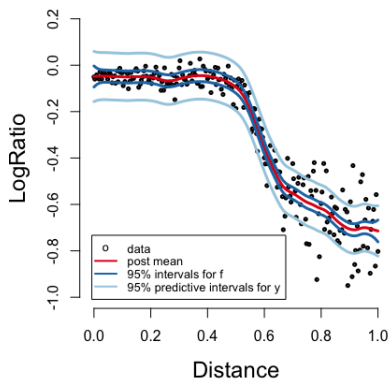
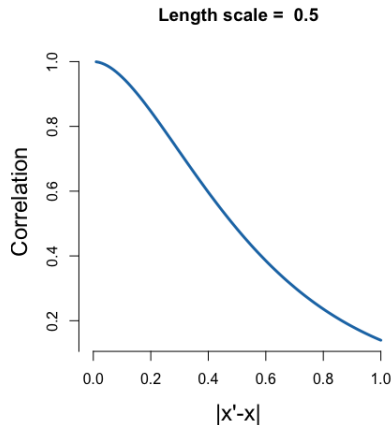
- Complete proof by result:

$$x_2 | x_1 = x_1^* \sim N [\mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (x_1^* - \mu_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}]$$

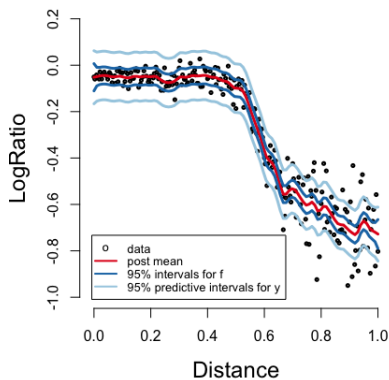
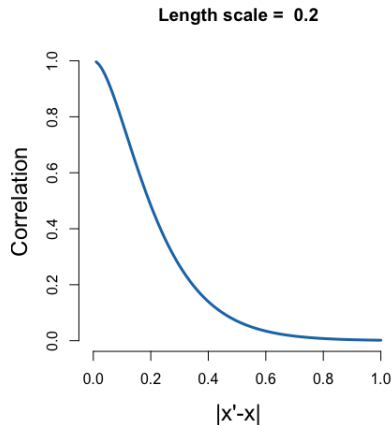
Example - LIDAR data



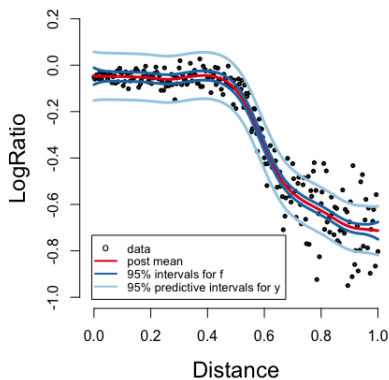
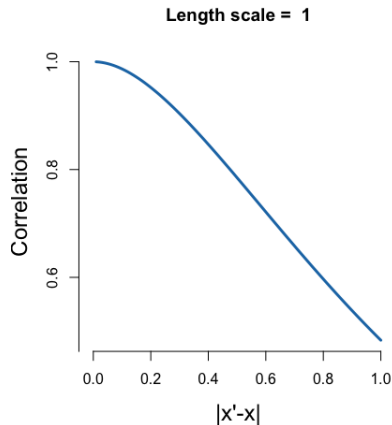
GP fit to LIDAR data $\ell = 0.5, \sigma_f = 0.5, \sigma_n = 0.05$



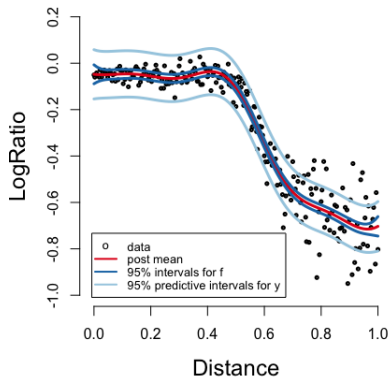
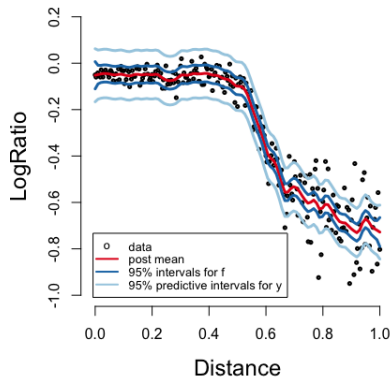
GP fit to LIDAR data $\ell = 0.2, \sigma_f = 0.5, \sigma_n = 0.05$



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

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

- Common: maximize the **marginal likelihood** wrt θ :

$$p(\mathbf{y}|\mathbf{X}, \theta) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{f}, \theta) p(\mathbf{f}|\mathbf{X}, \theta) d\mathbf{f}$$

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

- For **GP regression**: $\mathbf{y}|\mathbf{X}, \theta \sim N(0, K + \sigma_n^2 I)$ so

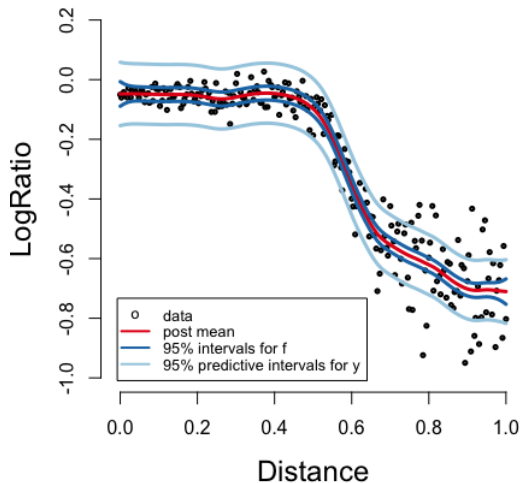
$$\log p(\mathbf{y}|\mathbf{X}, \theta) = -\frac{1}{2} \mathbf{y}^T (K + \sigma_n^2 I)^{-1} \mathbf{y} - \frac{1}{2} \log |K + \sigma_n^2 I| - \frac{n}{2} \log(2\pi)$$

- Proper **Bayesian inference** for hyperparameters (HMC?)

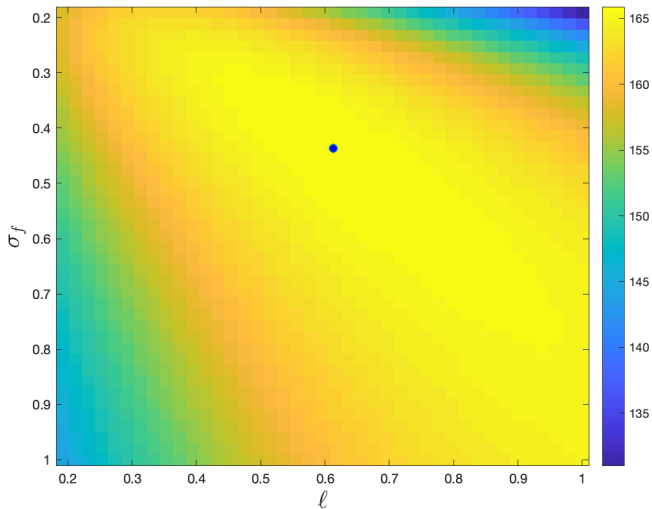
$$p(\theta|\mathbf{y}, \mathbf{X}) \propto p(\mathbf{y}|\mathbf{X}, \theta) p(\theta).$$

- Choice of kernel family by Bayesian model inference. For kernel $K_i \in \mathcal{K}$: $p(K_i|\mathbf{y}, \mathbf{X}) \propto p(\mathbf{y}|\mathbf{X}, K_i) p(K_i)$.

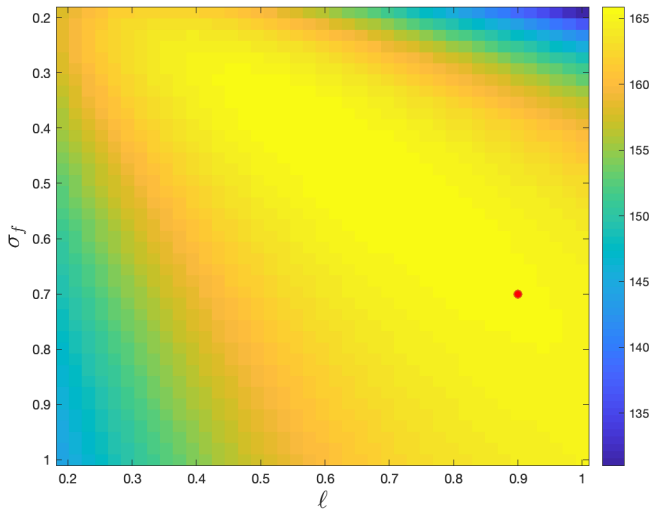
GP fit LIDAR $\ell_{opt} = 0.61, \sigma_{f,opt} = 0.44, \sigma_n = 0.05$



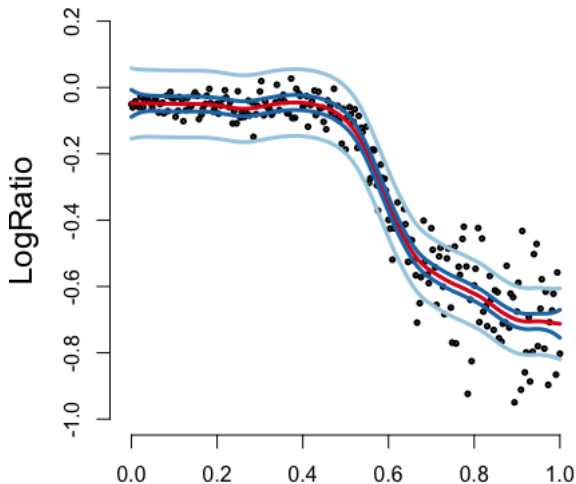
log marginal likelihood surface $\sigma_n = 0.05$



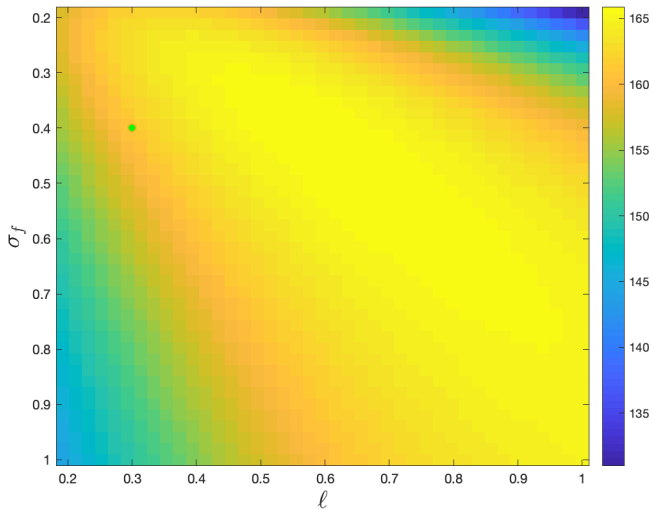
log marginal likelihood surface $\sigma_n = 0.05$



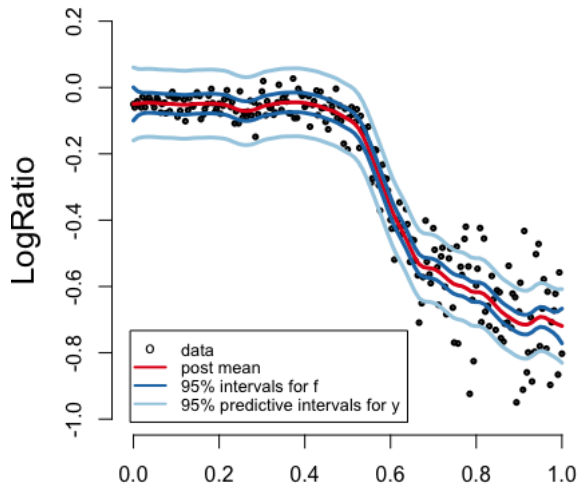
GP fit to LIDAR data $\ell = 0.9, \sigma_f = 0.70, \sigma_n = 0.05$



log marginal likelihood surface $\sigma_n = 0.05$



GP fit to LIDAR data $\ell = 0.3, \sigma_f = 0.4, \sigma_n = 0.05$



GP computations

- Covariance matrix K often numerically singular.
- Noise helps: $K + \sigma_n^2 I$.
- Artificial **jittering** $K + \epsilon I$ for small ϵ .
- Algorithm 2.1 and 3.1 in GPML for stable computations.
- We need to compute:
 - ▶ $(K + \sigma_n^2 I)^{-1} \mathbf{y}$ (posterior)
 - ▶ $\mathbf{y}^T (K + \sigma_n^2 I)^{-1} \mathbf{y}$ (log marginal likelihood)
 - ▶ $|\log K + \sigma_n^2 I|$ (log marginal likelihood)
- $(K + \sigma_n^2 I)^{-1} \mathbf{y}$ corresponds to solving $(K + \sigma_n^2 I) \mathbf{x} = \mathbf{y}$ wrt \mathbf{x} .

Cholesky + Forward and Backward substitution

- Solving $\mathbf{Ax} = \mathbf{b}$ by $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ is numerically unstable.
- **Cholesky factorization** $\mathbf{A} = \mathbf{LL}^T$ where \mathbf{L} is lower triangular.
- Let $\mathbf{Ax} = \mathbf{LL}^T\mathbf{x} = \mathbf{Lz} = \mathbf{b}$ if we define $\mathbf{z} = \mathbf{L}^T\mathbf{x}$.
- Solve $\mathbf{Lz} = \mathbf{b}$ wrt \mathbf{z} by forward substitution: $\mathbf{z} = \mathbf{L} \backslash \mathbf{b}$
- Solve $\mathbf{L}^T\mathbf{x} = \mathbf{z}$ wrt \mathbf{x} by backward substitution: $\mathbf{x} = \mathbf{L}^T \backslash \mathbf{z}$.
- So: $\mathbf{x} = \mathbf{L}^T \backslash (\mathbf{L} \backslash \mathbf{b})$.
- $|\mathbf{A}| = |\mathbf{LL}^T| = |\mathbf{L}|^2 = (\prod_{i=1}^p L_{ii})^2$.
- $\mathbf{y}^T \mathbf{A}^{-1} \mathbf{y} = \mathbf{y}^T (\mathbf{LL}^T)^{-1} \mathbf{y} = (\mathbf{L} \backslash \mathbf{y})^T (\mathbf{L} \backslash \mathbf{y})$.
- Cholesky also preserves sparsity.¹
- **Pre-conditioned conjugate gradient (PCG)**. $\mathbf{Ax} \approx \mathbf{b}$. Fast.

¹Rue and Held (2005). Gaussian Markov random fields. C&H.