

BAYESIAN LEARNING - LECTURE 6

Mattias Villani

**Division of Statistics
Department of Computer and Information Science
Linköping University**

LECTURE OVERVIEW

- ▶ Flexible nonlinear regression and splines
- ▶ Smoothness/shrinkage priors
- ▶ Gaussian process regression

POLYNOMIAL REGRESSION

- Recall the linear regression model with a single covariate

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2).$$

- Polynomial regression:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_k x_i^k + \varepsilon_i.$$

This can be written as a linear regression

$$y = X_P \beta + \varepsilon,$$

where

$$X_P = (1, x, x^2, \dots, x^k).$$

- The posterior of β is obtained like any linear regression.

SPLINES

- ▶ A more local basis: **truncated polynomials**

$$b_{ij} = \begin{cases} (x_i - k_j)^p & \text{if } x_i > k_j \\ 0 & \text{otherwise} \end{cases}$$

- ▶ k_1, k_2, \dots, k_m are the **knots**.
- ▶ Splines are nonlinear in x , but linear in basis space

$$y = X_b \beta + \varepsilon,$$

where X_b is the basis regression matrix

$$X_b = (b_1, \dots, b_m).$$

- ▶ Common extension

$$X_b = (1, x, b_1, \dots, b_m).$$

- ▶ Still just linear in X_b . Linear regression fitting.

SMOOTHNESS PRIOR FOR SPLINES

- ▶ Problem: too many knots leads to **over-fitting**.
- ▶ Solution: **smoothness/shrinkage/regularization prior**

$$\beta_i | \sigma^2 \stackrel{iid}{\sim} N\left(0, \frac{\sigma^2}{\lambda}\right)$$

- ▶ Larger λ gives smoother fit. Note: here we have $\Omega_0 = \lambda I$.
- ▶ Equivalent to a penalized likelihood:

$$-2 \cdot \text{LogPost} \propto \text{RSS}(\beta) + \lambda \beta' \beta$$

- ▶ Posterior mean gives **ridge regression** estimator

$$\tilde{\beta} = (X'X + \lambda I)^{-1} X'y$$

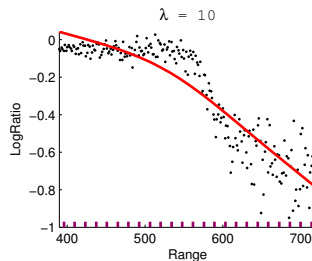
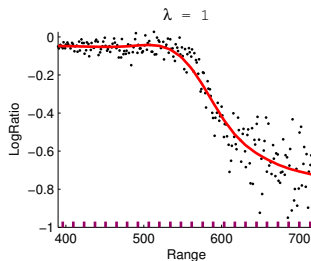
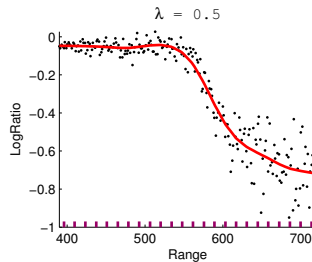
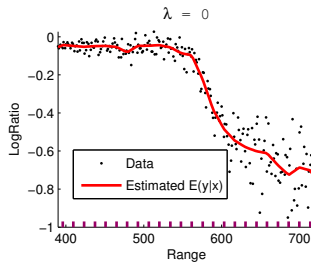
- ▶ **Shrinkage** toward zero

$$\text{As } \lambda \rightarrow \infty, \tilde{\beta} \rightarrow 0$$

- ▶ When $X'X = I$

$$\tilde{\beta} = \frac{1}{1 + \lambda} \hat{\beta}_{OLS}$$

BAYESIAN SPLINE WITH SMOOTHNESS PRIOR



SMOOTHNESS PRIOR FOR SPLINES, CONT.

- ▶ The famous **Lasso** variable selection method is equivalent to using the posterior mode estimate under the prior:

$$\beta_i | \sigma^2 \stackrel{iid}{\sim} \text{Laplace} \left(0, \frac{\sigma^2}{\lambda} \right)$$

where the general Laplace density is

$$p(\beta_i) = \frac{1}{2b} \exp \left(-\frac{|\beta_i - \mu|}{b} \right)$$

- ▶ The Bayesian shrinkage prior is **interpretable, not ad hoc**.
- ▶ Laplace distribution have heavy tails.
- ▶ Laplace prior: many β_i close to zero, but some β_i may be very large.
- ▶ Normal distribution have light tails.
- ▶ Normal prior: most β_i are fairly equal in size, and no single β_i can be very much larger than the other ones.

ESTIMATING THE SHRINKAGE

- ▶ How do we determine the degree of smoothness, λ ? Cross-validation.
- ▶ Bayesian: I cannot specify $\lambda \Rightarrow \lambda$ is unknown \Rightarrow use a prior for λ .
- ▶ One possibility: $\lambda \sim \text{Inv} - \chi^2(\eta_0, \lambda_0)$. The user specifies η_0 and λ_0 .
- ▶ Alternative approach: specify the prior on the *degrees of freedom*.
- ▶ Hierarchical setup:

$$\begin{aligned}y|\beta, x &\sim N(x'\beta, \sigma^2) \\ \beta|\sigma^2 &\sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \begin{pmatrix} \delta_0^{-1} I_q & 0 \\ 0 & \lambda^{-1} I_m \end{pmatrix}\right) \\ \sigma^2 &\sim \text{Inv} - \chi^2(\nu_0, \sigma_0^2) \\ \lambda &\sim \text{Inv} - \chi^2(\eta_0, \lambda_0)\end{aligned}$$

Note: different shrinkage on the original q covariates (δ_0) and the covariates that comes from the knots (λ).

ESTIMATING THE SHRINKAGE, CONT.

- ▶ Joint posterior

$$p(\beta, \sigma^2, \lambda | y, x) = p(\beta, \sigma^2 | \lambda, y, x) p(\lambda | y, x)$$

where

$$p(\lambda | y, x) = \int \int p(\beta, \sigma, \lambda | y, x) d\beta d\sigma^2$$

is the marginal posterior of λ .

- ▶ The conditional posterior $p(\beta, \sigma^2 | \lambda, y, x)$ is a special case linear regression with conjugate prior $\mu_0 = (0, 0)'$ and

$$\Omega_0 = \begin{pmatrix} \delta_0 I_q & 0 \\ 0 & \lambda I_m \end{pmatrix}$$

SUMMARY OF THE POSTERIOR WITH NORMAL SHRINKAGE PRIOR

- The joint posterior of β , σ^2 and λ is

$$\beta|\sigma^2, \lambda, y \sim N(\mu_n, \Omega_n^{-1})$$

$$\sigma^2|\lambda, y \sim \text{Inv} - \chi^2(\nu_n, \sigma_n^2)$$

$$p(\lambda|y) \propto \sqrt{\frac{|\Omega_0|}{|X'X + \Omega_0|}} \left(\frac{\nu_n \sigma_n^2}{2}\right)^{-\nu_n/2} \cdot p(\lambda)$$

where $p(\lambda)$ is the prior for λ , and

$$\mu_n = (X'X + \Omega_0)^{-1} X'y$$

$$\Omega_n = X'X + \Omega_0$$

$$\nu_n = \nu_0 + n$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + y'y - \mu_n' \Omega_n \mu_n$$

BAYESIAN VARIABLE SELECTION AND ESTIMATING KNOT LOCATIONS

- ▶ Selecting among a set of fixed knots k_1, \dots, k_m is a variable selection problem. More on **Bayesian variable selection** in the last module.
- ▶ The **location of the knots** can be treated as **unknown**, and estimated from the data.
- ▶ The joint posterior of parameters and knot locations

$$p(\beta, \sigma^2, \lambda, \xi_1, \dots, \xi_q | y, x)$$

where ξ_i is the location of the i th knot.

- ▶ Posterior is complex but can be sampled from by Markov Chain Monte Carlo (MCMC). Li and Villani (2013, SJS).

NONPARAMETRIC REGRESSION - SMOOTH INTERPOLATION

- ▶ Another approach treats all n ordinates as unknown parameters:

$$f(x_i) = \gamma_i.$$

- ▶ Problem: too many parameters. Estimated curve wiggles way too much.
- ▶ Solution: use a (multivariate) prior on $\gamma = (\gamma_1, \dots, \gamma_n)'$ that carries the info that the regression curve is smooth:

if x_i and x_k are close then γ_i is close to γ_k

- ▶ Order the data with respect to covariates and assign the prior

$$p(\gamma_i | \gamma_{i-1}) \sim N(\gamma_{i-1}, \tau_0^2 \cdot |x_i - x_{i-1}|), \text{ for } i = 2, \dots, n.$$

- ▶ The hyperparameter τ_0^2 controls the degree of prior smoothness.

NON-PARAMETRIC REGRESSION

- ▶ **Linear regression**

$$y = \beta \cdot x + \varepsilon$$

where $\varepsilon \sim N(0, \sigma^2)$ and iid over observations.

- ▶ **Nonlinear regression**

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

where $f(\cdot)$ is some nonlinear function (ex $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$).

NON-PARAMETRIC REGRESSION

- ▶ **Linear regression**

$$y = \beta \cdot x + \varepsilon$$

where $\varepsilon \sim N(0, \sigma^2)$ and iid over observations.

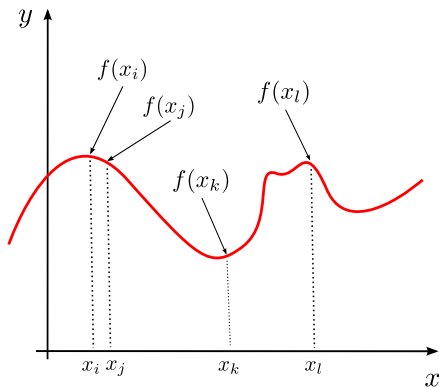
- ▶ **Nonlinear regression**

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

where $f(\cdot)$ is some nonlinear function (ex $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$).

- ▶ **Non-parametric regression:** avoiding a parametric form for $f(\cdot)$.
- ▶ How do we put a **prior over a set of functions**?
- ▶ Restrict attention to a grid of (ordered) x -values: x_1, x_2, \dots, x_k . [weight space view].
- ▶ We can now put a joint prior on the k function values:
 $f(x_1), f(x_2), \dots, f(x_k)$.

NONPARAMETRIC = ONE PARAMETER FOR EVERY x !



GAUSSIAN PROCESS REGRESSION

- ▶ Natural choice. Multivariate normal (Gaussian):

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

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

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

GAUSSIAN PROCESS REGRESSION

- Natural choice. Multivariate normal (Gaussian):

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

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

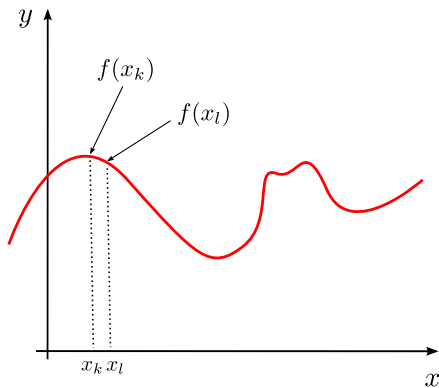
$$\text{Cov}(f(x_p), f(x_q))$$

- **Squared exponential covariance function**

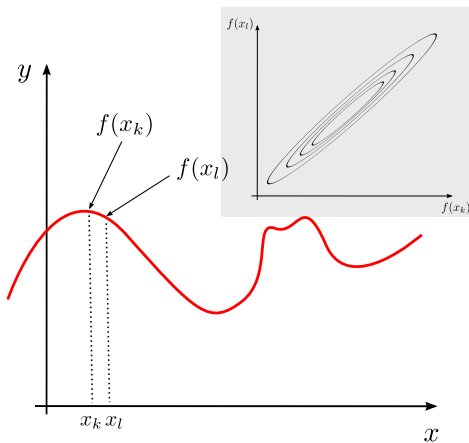
$$\text{Cov}(f(x_p), f(x_q)) = K(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2} \left(\frac{x_p - x_q}{\ell}\right)^2\right)$$

- The covariance between $f(x_p)$ and $f(x_q)$ is a function of x_p and x_q .
- 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 (no need for a pre-determined grid)

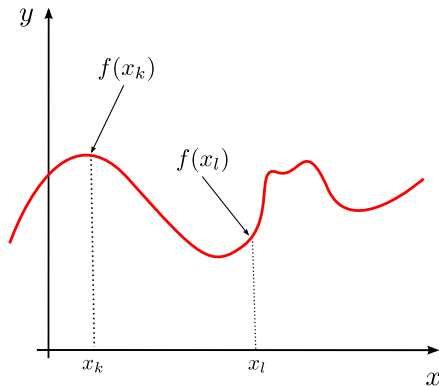
SMOOTH FUNCTION - POINTS NEARBY



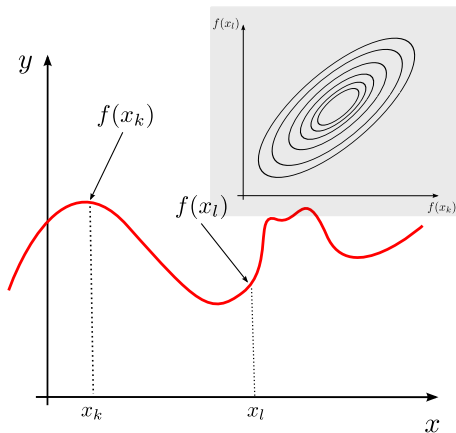
SMOOTH FUNCTION - POINTS NEARBY



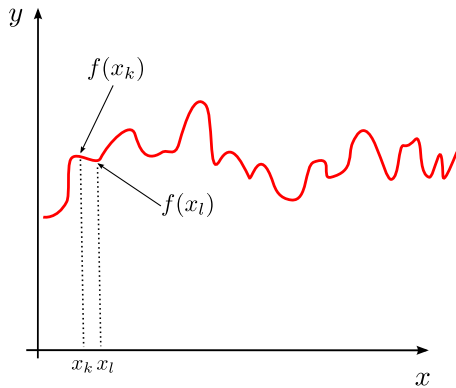
SMOOTH FUNCTION - POINTS FAR APART



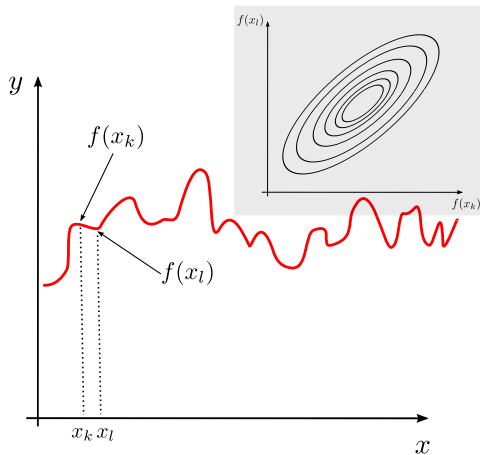
SMOOTH FUNCTION - POINTS FAR APART



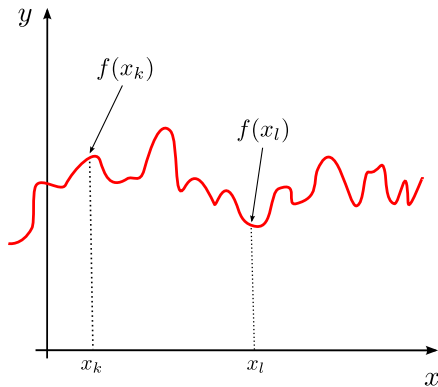
JAGGED FUNCTION - POINTS NEARBY



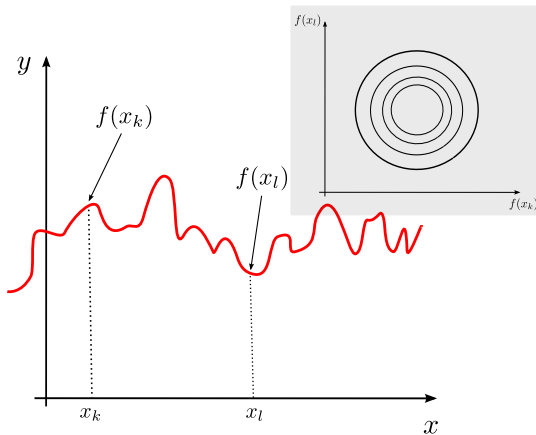
JAGGED FUNCTION - POINTS NEARBY



JAGGED FUNCTION - POINTS FAR APART



JAGGED FUNCTION - POINTS FAR APART



GAUSSIAN PROCESS REGRESSION, CONT.

DEFINITION

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

- ▶ A Gaussian process is really a **probability distribution over functions** (curves). This is exactly what we want! No need for a grid!
- ▶ A GP is completely 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' (note: this is *not* the transpose here).

- ▶ A Gaussian process (prior) is denoted by

$$f(x) \sim GP(m(x), K(x, x'))$$

GAUSSIAN PROCESS REGRESSION, CONT.

► Example:

$$m(x) = \sin(x)$$

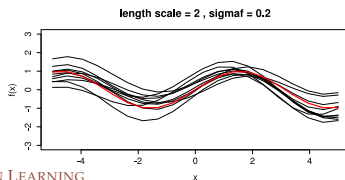
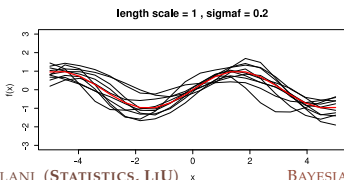
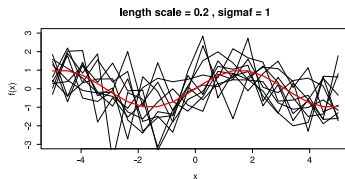
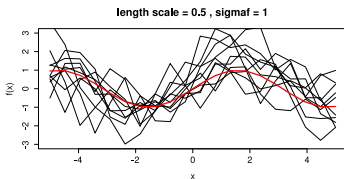
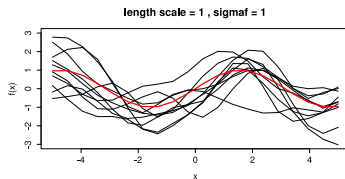
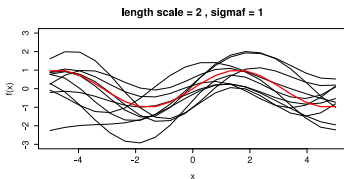
$$K(x, x') = \sigma_f^2 \exp \left(-\frac{1}{2} \left(\frac{x_p - x_q}{\ell} \right)^2 \right)$$

where $\ell > 0$ is the length scale.

- Larger ℓ gives more smoothness in $f(x)$.
- Simulate draw from $f(x) \sim GP(m(x), K(x, x'))$ over any grid $x_* = (x_1, \dots, x_n)$ by using that

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

SIMULATING A GP - SINE MEAN AND SE KERNEL



GAUSSIAN PROCESS REGRESSION, CONT.

► Model

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma^2)$$

► Prior

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

- You have observed the data: $\mathbf{x} = (x_1, \dots, x_n)'$ and $\mathbf{y} = (y_1, \dots, y_n)'$.
- Goal: the posterior of $f(\cdot)$ over a grid of x -values: $\mathbf{f}_* = \mathbf{f}(\mathbf{x}_*)$.

GAUSSIAN PROCESS REGRESSION, CONT.

► Model

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma^2)$$

► Prior

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

► You have observed the data: $\mathbf{x} = (x_1, \dots, x_n)'$ and $\mathbf{y} = (y_1, \dots, y_n)'$.

► Goal: the posterior of $f(\cdot)$ over a grid of x -values: $\mathbf{f}_* = \mathbf{f}(\mathbf{x}_*)$.

► Intermediate step: joint distribution of \mathbf{y} and \mathbf{f}_*

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim N \left\{ \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{bmatrix} K(\mathbf{x}, \mathbf{x}) + \sigma^2 I & K(\mathbf{x}, \mathbf{x}_*) \\ K(\mathbf{x}_*, \mathbf{x}) & K(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right\}$$

GAUSSIAN PROCESS REGRESSION, CONT.

► Model

$$y_i = f(\mathbf{x}_i) + \varepsilon_i, \quad \varepsilon \stackrel{iid}{\sim} N(0, \sigma^2)$$

► Prior

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

► You have observed the data: $\mathbf{x} = (x_1, \dots, x_n)'$ and $\mathbf{y} = (y_1, \dots, y_n)'$.

► Goal: the posterior of $f(\cdot)$ over a grid of x -values: $\mathbf{f}_* = \mathbf{f}(\mathbf{x}_*)$.

► Intermediate step: joint distribution of \mathbf{y} and \mathbf{f}_*

$$\begin{pmatrix} \mathbf{y} \\ \mathbf{f}_* \end{pmatrix} \sim N \left\{ \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{bmatrix} K(\mathbf{x}, \mathbf{x}) + \sigma^2 I & K(\mathbf{x}, \mathbf{x}_*) \\ K(\mathbf{x}_*, \mathbf{x}) & K(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right\}$$

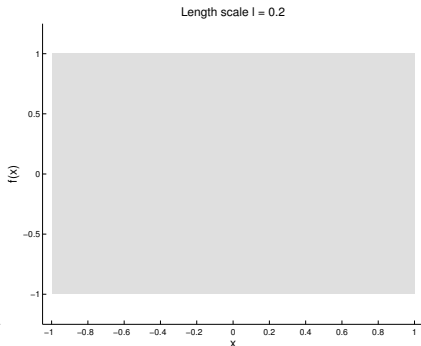
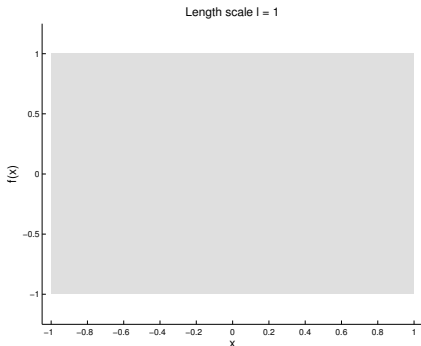
► The posterior

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

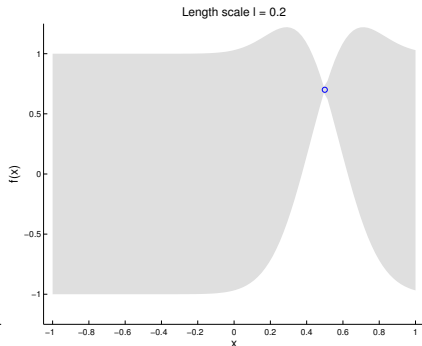
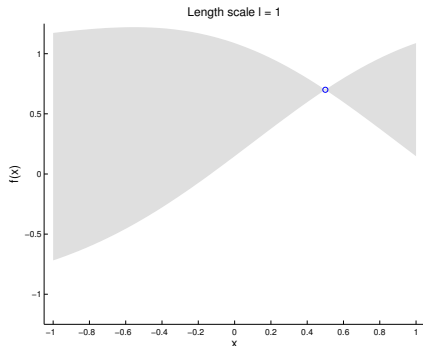
$$\bar{\mathbf{f}}_* = K(\mathbf{x}_*, \mathbf{x}) [K(\mathbf{x}, \mathbf{x}) + \sigma^2 I]^{-1} \mathbf{y}$$

$$\text{cov}(\mathbf{f}_*) = K(\mathbf{x}_*, \mathbf{x}_*) - K(\mathbf{x}_*, \mathbf{x}) [K(\mathbf{x}, \mathbf{x}) + \sigma^2 I]^{-1} K(\mathbf{x}, \mathbf{x}_*)$$

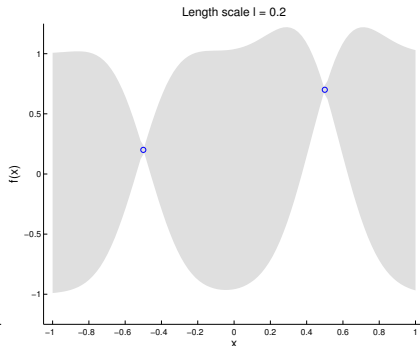
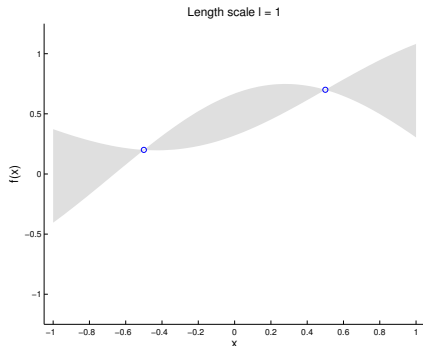
LEARNING A NOISE-FREE GAUSSIAN PROCESS



LEARNING A NOISE-FREE GAUSSIAN PROCESS



LEARNING A NOISE-FREE GAUSSIAN PROCESS



LEARNING A NOISE-FREE GAUSSIAN PROCESS

