

INTRODUCTION TO MACHINE LEARNING

TOPIC 6: GAUSSIAN PROCESSES AND MIXTURE MODELS

LECTURE 1

Mattias Villani

Division of Statistics and Machine Learning
Department of Computer and Information Science
Linköping University



TOPIC OVERVIEW

- ▶ Recall: **The multivariate normal distribution**
- ▶ Bayesian inference for **Gaussian linear/nonlinear regression**.
- ▶ **Gaussian processes** for **nonparametric regression**
 - ▶ Covariance kernels
 - ▶ Selecting the kernel and hyperparameters
- ▶ **Introduction to GP classification**

THE MULTIVARIATE NORMAL DISTRIBUTION

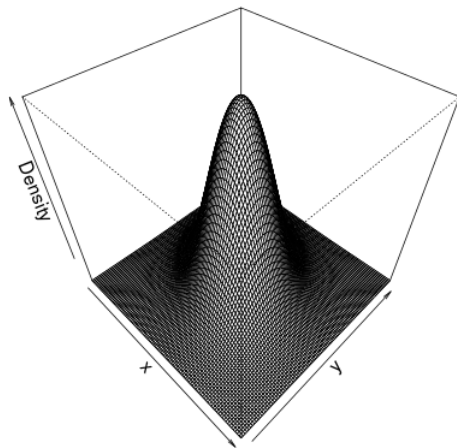
- The **density function** of a p -variate normal vector $\mathbf{x} \sim N(\mu, \Sigma)$ is

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

- Example: **Bivariate normal** ($p = 2$)

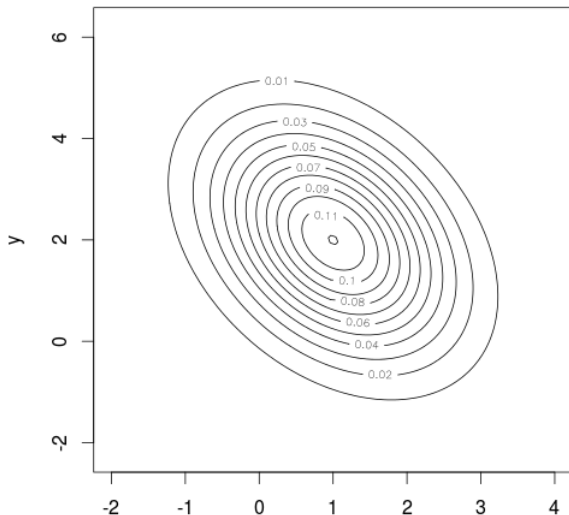
$$\Sigma = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}$$

MULTIVARIATE NORMAL

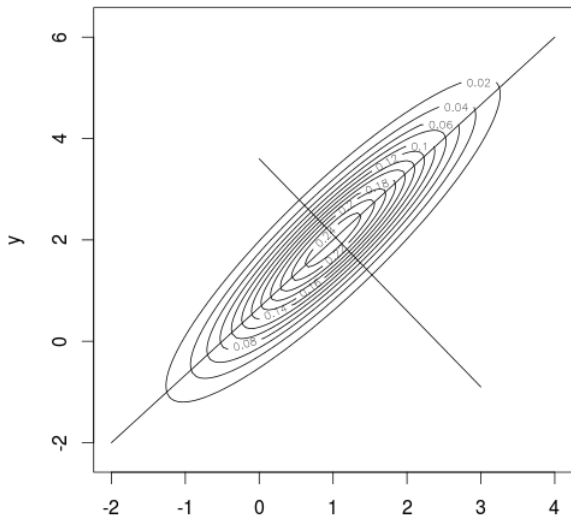


MULTIVARIATE NORMAL

Marginals are normal, joint is normal



MULTIVARIATE NORMAL



FLEXIBLE NONLINEAR REGRESSION

- ▶ **Linear regression**

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

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

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

- ▶ The weights \mathbf{w} are called regression coefficients (β) in statistics.
- ▶ **Polynomial regression**: $\phi(\mathbf{x}) = (1, x, x^2, x^3, \dots, x^k)$:

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

- ▶ More generally: **splines** with **basis functions**. See Topic 7 in this course.

BAYESIAN LINEAR REGRESSION - INFERENCE

- ▶ Linear regression for all n observations

$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\mathbf{w}} + \underset{n \times 1}{\boldsymbol{\varepsilon}}$$

- ▶ \mathbf{w} is unknown. σ_n is assumed known.

- ▶ **Prior**

$$\mathbf{w} \sim N(0, \Sigma_p)$$

- ▶ Ridge regression: $\Sigma_p = \lambda^{-1}I$.

- ▶ **Posterior**

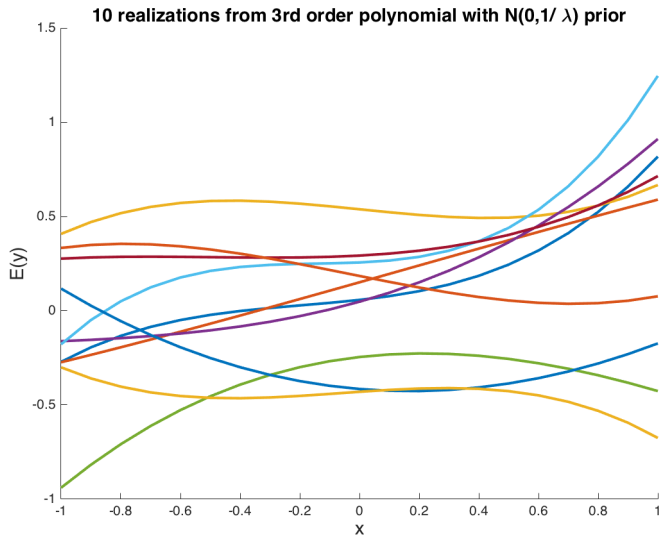
$$\mathbf{w}|\mathbf{X}, \mathbf{y} \sim N(\bar{\mathbf{w}}, \mathbf{A}^{-1})$$

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

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

- ▶ Recall: **Posterior precision = Data Precision + Prior Precision.**

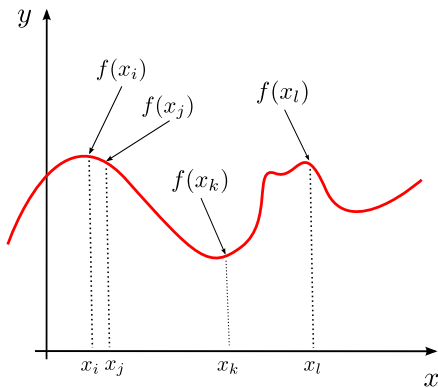
FLEXIBLE NONLINEAR REGRESSION



NON-PARAMETRIC REGRESSION

- ▶ **Non-parametric regression:** avoiding a parametric form for $f(\cdot)$.
Treat $f(\mathbf{x})$ as an unknown parameter for every \mathbf{x} .
- ▶ **Weight space view**
 - ▶ Restrict attention to a grid of (ordered) x -values: x_1, x_2, \dots, x_k .
 - ▶ Put a joint prior on the k function values: $f(x_1), f(x_2), \dots, f(x_k)$.
- ▶ **Function space view**
 - ▶ Treat f as an **unknown function**.
 - ▶ Put a **prior over a set of functions**.

NONPARAMETRIC = ONE PARAMETER FOR EVERY x !



GAUSSIAN PROCESS REGRESSION

- ▶ Weight-space view. GP assumes

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

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

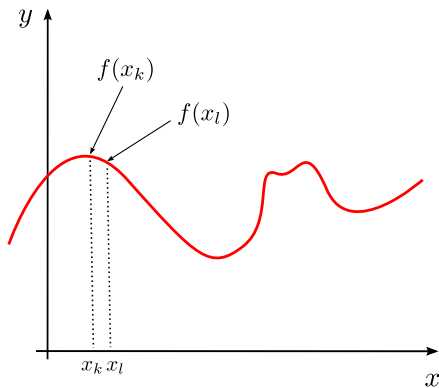
$$\text{Cov}(f(x_p), f(x_q))$$

- ▶ **Squared exponential** covariance function

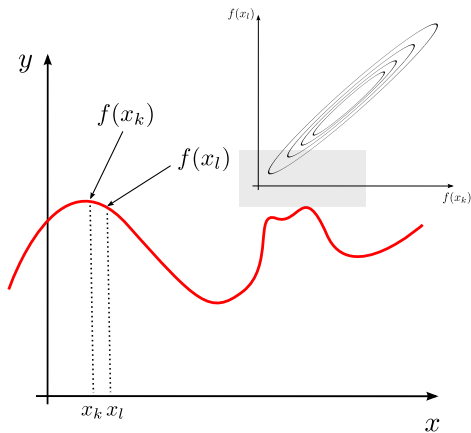
$$\text{Cov}(f(x_p), f(x_q)) = k(x_p, x_q) = \sigma_f^2 \exp\left(-\frac{1}{2}(x_p - x_q)^2\right)$$

- ▶ 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 .
- ▶ Extension to multiple covariates: $(x_p - x_q)^2$ replaced by $(\mathbf{x}_p - \mathbf{x}'_q)^T(\mathbf{x}_p - \mathbf{x}'_q)$.

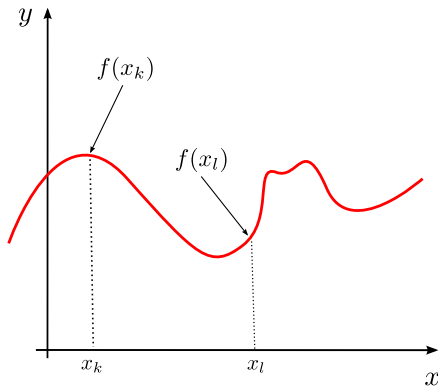
SMOOTH FUNCTION - POINTS NEARBY



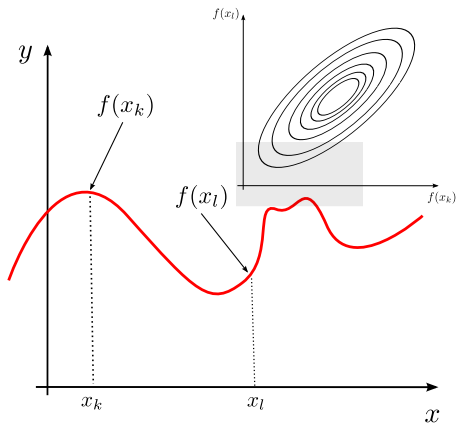
SMOOTH FUNCTION - POINTS NEARBY



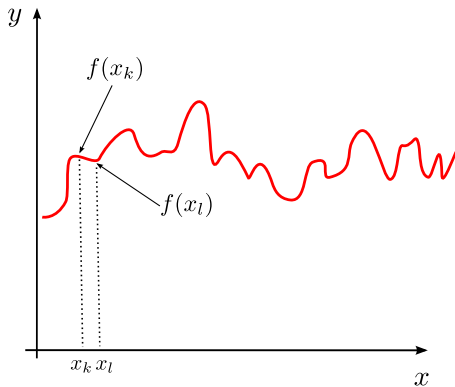
SMOOTH FUNCTION - POINTS FAR APART



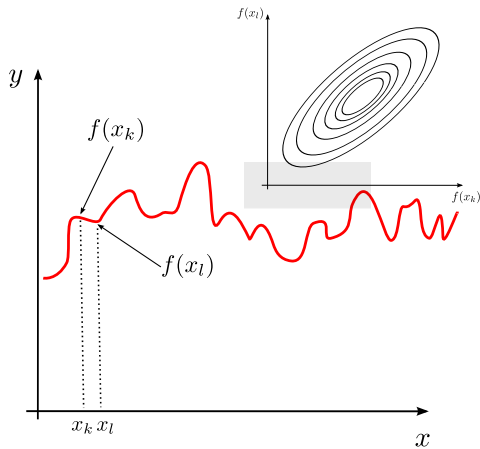
SMOOTH FUNCTION - POINTS FAR APART



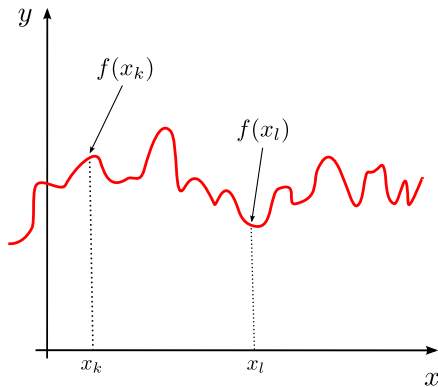
JAGGED FUNCTION - POINTS NEARBY



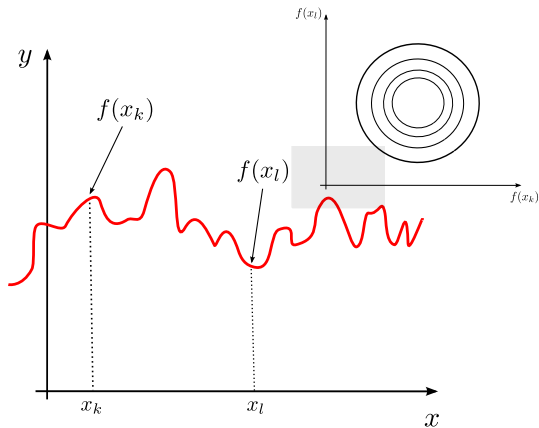
JAGGED FUNCTION - POINTS NEARBY



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).
- ▶ 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** is denoted by

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

- ▶ **Bayesian:** $f(x) \sim GP$ encodes **prior beliefs** about the unknown $f(\cdot)$.

A SIMPLE GP EXAMPLE

- ▶ Example:

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

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

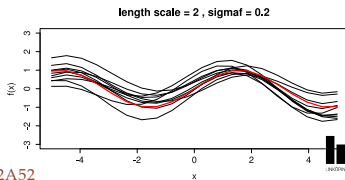
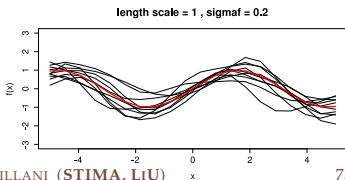
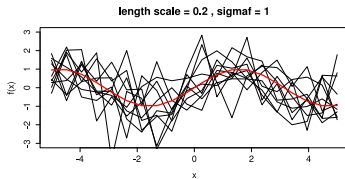
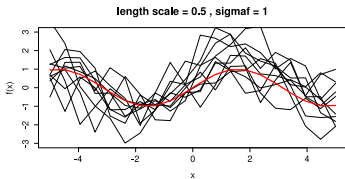
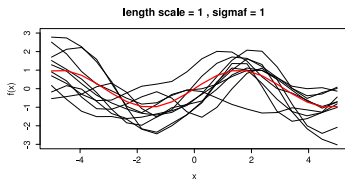
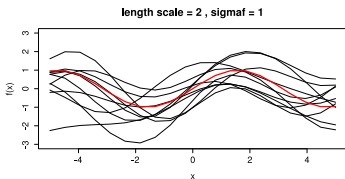
where $\ell > 0$ is the length scale.

- ▶ Larger ℓ gives more smoothness in $f(x)$.
- ▶ Simulate draw from $f(x) \sim GP(m(x), k(x, x'))$ over a grid $\mathbf{x}_* = (x_1, \dots, x_n)$ by using that

$$f(\mathbf{x}_*) \sim N(m(\mathbf{x}_*), K(\mathbf{x}_*, \mathbf{x}_*))$$

- ▶ Note that the **kernel** $k(x, x')$ produces a **covariance matrix** $K(\mathbf{x}_*, \mathbf{x}_*)$ when evaluated at the vector \mathbf{x}_* .

SIMULATING A GP - SINE MEAN AND SE KERNEL



SIMULATING A GP

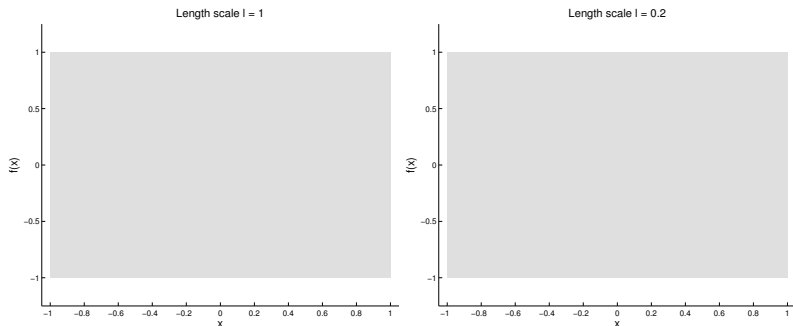
- The joint way: Choose a grid x_1, \dots, x_k . Simulate the k -vector

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

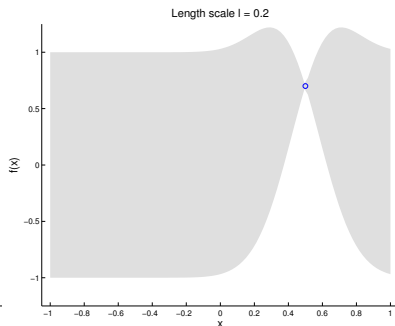
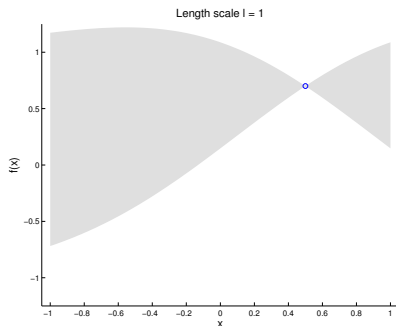
- More intuition from the conditional decomposition

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

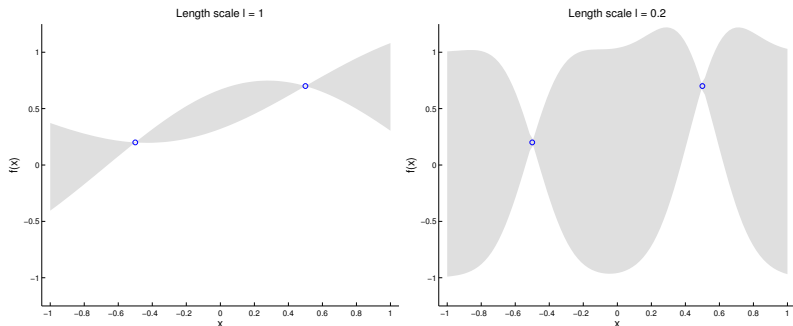
SIMULATION FROM $\ell=1$ VS $\ell=0.2$. BEFORE FIRST DRAW.



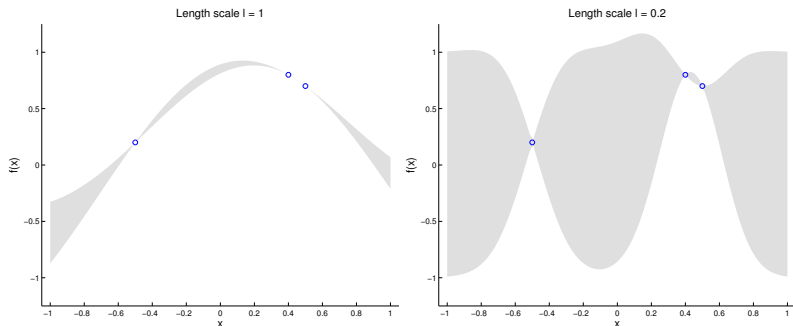
SIMULATION FROM $\ell=1$ VS $\ell=0.2$. BEFORE SECOND DRAW.



SIMULATION FROM $\ell=1$ VS $\ell=0.2$. BEFORE THIRD DRAW.



SIMULATION FROM $\ell=1$ VS $\ell=0.2$. BEFORE FOURTH DRAW.



THE POSTERIOR FOR A GAUSSIAN PROCESS REGRESSION

► Model

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

► Prior

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

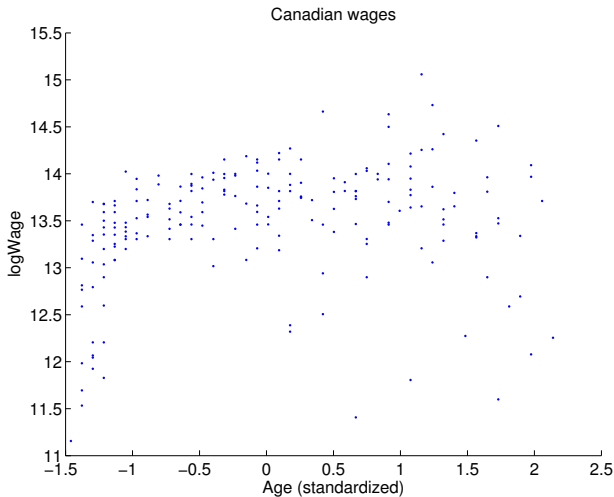
- You have observed the data: $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$ and $\mathbf{y} = (y_1, \dots, y_n)'$.
- Goal: the posterior of $f(\cdot)$ over a grid of \mathbf{x} -values: $\mathbf{f}_* = \mathbf{f}(\mathbf{x}_*)$.
- The **posterior** (use formula for conditional Gaussian above)

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

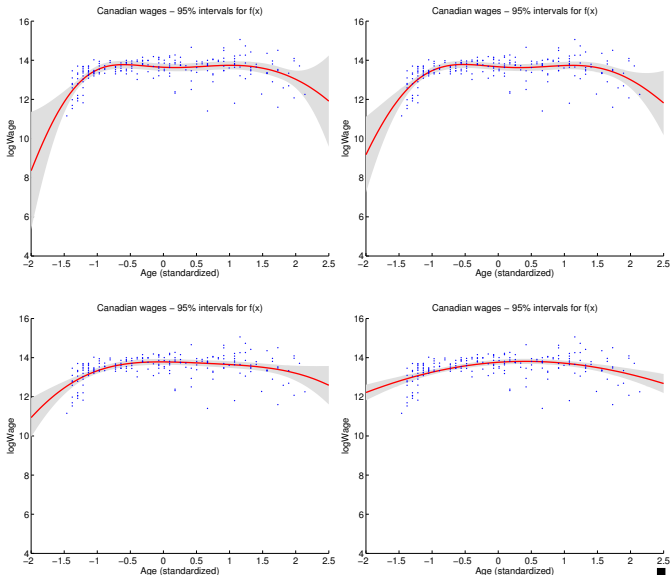
$$\bar{\mathbf{f}}_* = K(\mathbf{x}_*, \mathbf{x}) [K(\mathbf{x}, \mathbf{x}) + \sigma^2 I]^{-1} \mathbf{y}$$

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

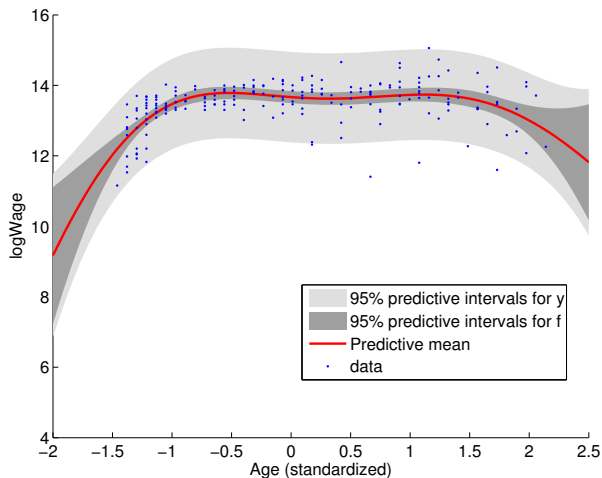
EXAMPLE - CANADIAN WAGES



POSTERIOR OF $F - \ell = 0.2, 0.5, 1, 2$



CANADIAN WAGES - PREDICTION WITH $\ell = 0.5$



TWO COMMONLY USED COVARIANCE KERNELS

- ▶ Let $r = \|x - x'\|$.
- ▶ **Squared exponential (SE)** ($\ell > 0, \sigma_f > 0$)

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

- ▶ Infinitely mean square differentiable. Very smooth.
- ▶ **Matérn** ($\ell > 0, \sigma_f > 0, \nu > 0$)

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

- ▶ $\nu = 3/2$ and $\nu = 5/2$ most useful for ML. As $\nu \rightarrow \infty$, Matérn's kernel approaches SE kernel.

MORE THAN ONE INPUT - ARD

- ▶ Anisotropic version of isotropic kernels by setting $r^2(\mathbf{x}, \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \mathbf{M} (\mathbf{x} - \mathbf{x}')$ where \mathbf{M} is positive definite.
- ▶ **Automatic Relevance Determination (ARD)**:
 $\mathbf{M} = \text{Diag}(\ell_1^{-2}, \dots, \ell_D^{-2})$ is diagonal with different length scales.
- ▶ ARD does 'variable selection' since large ℓ_j means that the j th input essentially drops out of $f(\mathbf{x})$.

DETERMINING THE HYPERPARAMETERS

- ▶ Kernel depends on **hyperparameters** θ . Example SE kernel $[\theta = (\sigma_f, \ell)^T]$

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

- ▶ Common approach: choose the hyperparameters that maximizes the **marginal likelihood** (**evidence**):

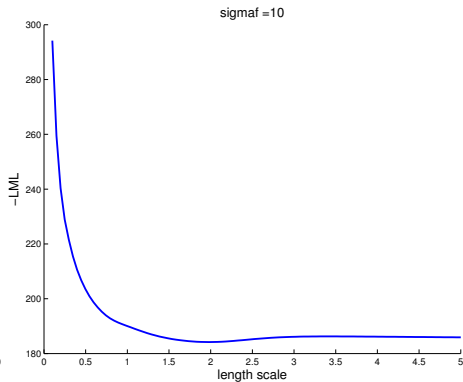
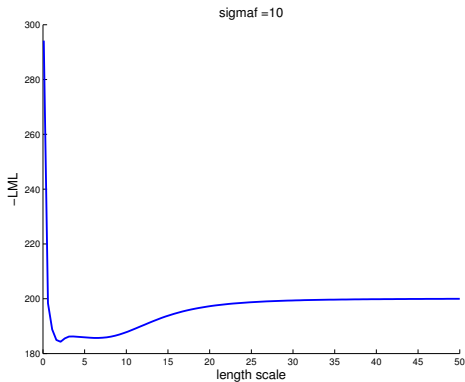
$$p(\mathbf{y}|\mathbf{X}, \theta) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{f}, \theta) p(\mathbf{f}|\mathbf{X}, \theta) d\mathbf{f}$$

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

- ▶ For Gaussian process regression:

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

CANADIAN WAGES - LML DETERMINATION OF ℓ



CLASSIFICATION WITH LOGISTIC REGRESSION

- ▶ **Classification:** binary response $y \in \{-1, 1\}$ predicted by features \mathbf{x} .
- ▶ Example: linear logistic regression

$$Pr(y = 1|\mathbf{x}) = \lambda(\mathbf{x}^T \mathbf{w})$$

where $\lambda(z)$ is the logistic **link function**

$$\lambda(z) = \frac{1}{1 + \exp(-z)}$$

- ▶ $\lambda(z)$ 'squashes' the linear prediction $\mathbf{x}^T \mathbf{w} \in \mathbb{R}$ into $\lambda(\mathbf{x}^T \mathbf{w}) \in [0, 1]$.
- ▶ Logistic regression has **linear decision boundaries**.

GP CLASSIFICATION

- ▶ Obvious **GP extension** of logistic regression: replace $\mathbf{x}^T \mathbf{w}$ by $f(\mathbf{x})$ where

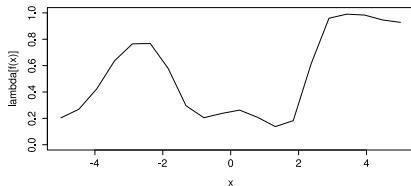
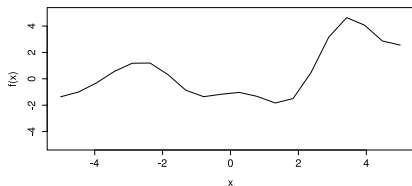
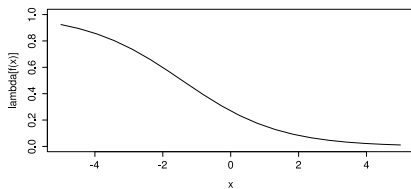
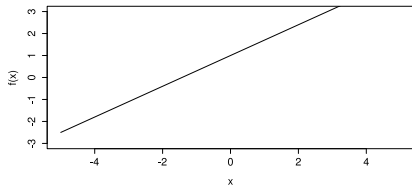
$$f(\mathbf{x}) \sim GP(0, k(\mathbf{x}, \mathbf{x}'))$$

and squash f through logistic function (or normal CDF)

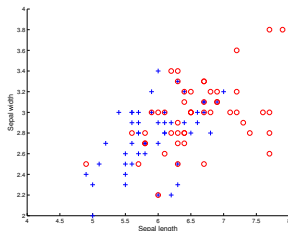
$$Pr(y = 1|\mathbf{x}) = \lambda(f(\mathbf{x}))$$

- ▶ Posterior and predictive distribution is complicated. Solutions:
 - ▶ Approximations: **Laplace**, **Expectation Propagation (EP)** or Variational Bayes (VB)
 - ▶ MCMC sampling.

SQUASHING F

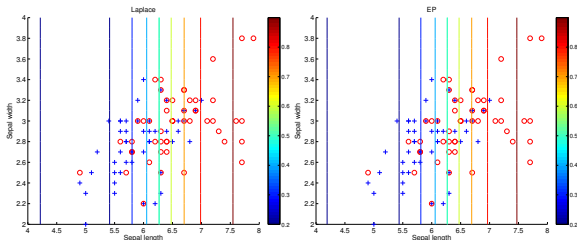


IRIS DATA - SEPAL - SE KERNEL WITH ARD

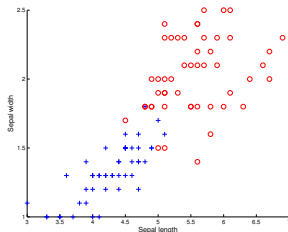


Laplace: $\hat{\ell}_1 = 1.7214, \hat{\ell}_2 = 185.5040, \sigma_f = 1.4361$

EP: $\hat{\ell}_1 = 1.7189, \hat{\ell}_2 = 55.5003, \sigma_f = 1.4343$

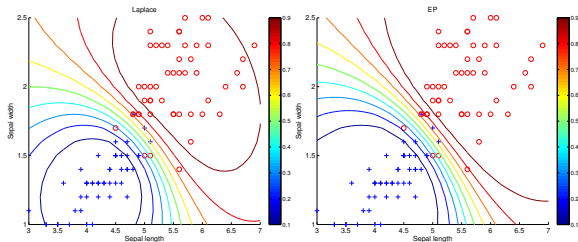


IRIS DATA - PETAL - SE KERNEL WITH ARD

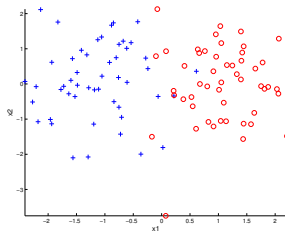


Laplace: $\hat{\ell}_1 = 1.7606, \hat{\ell}_2 = 0.8804, \sigma_f = 4.9129$

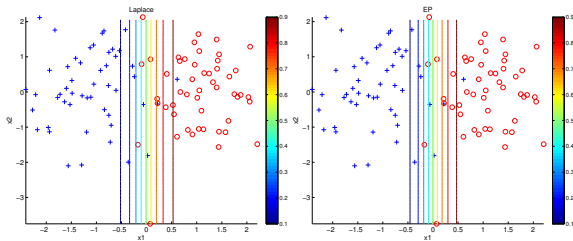
EP: $\hat{\ell}_1 = 2.1139, \hat{\ell}_2 = 1.0720, \sigma_f = 5.3369$



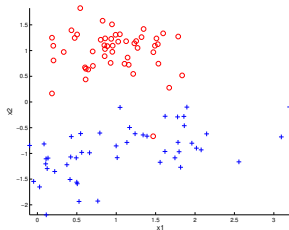
TOY DATA 1 - SE KERNEL WITH ARD



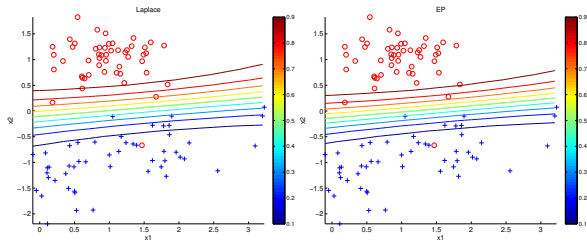
$$\text{EP: } \hat{\ell}_1 = 2.4503, \hat{\ell}_2 = 721.7405, \sigma_f = 4.7540$$



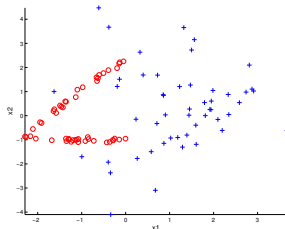
TOY DATA 2 - SE KERNEL WITH ARD



EP: $\hat{\ell}_1 = 8.3831, \hat{\ell}_2 = 1.9587, \sigma_f = 4.5483$



TOY DATA 3 - SE KERNEL WITH ARD



Laplace: $\hat{\ell}_1 = 0.7726, \hat{\ell}_2 = 0.6974, \sigma_f = 11.7854$

EP: $\hat{\ell}_1 = 1.2685, \hat{\ell}_2 = 1.0941, \sigma_f = 17.2774$

