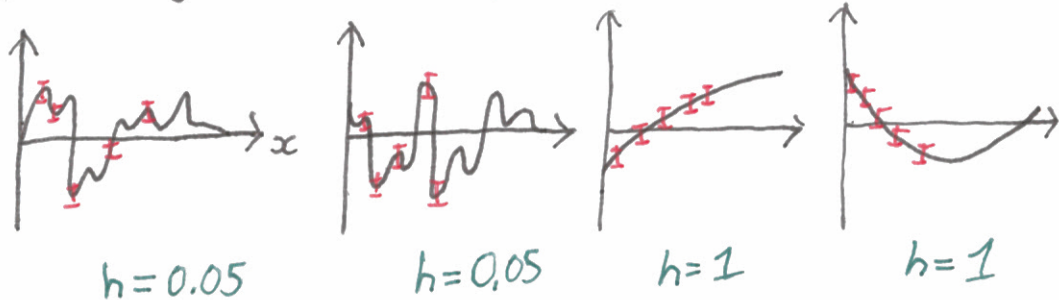


Two Bayesian Linear Regression models

Both $\underline{w} \sim N(\underline{0}, \sigma_w^2 \mathbb{I})$, $y \sim N(\Phi \underline{w}, \sigma_y^2)$
 from RBFs \uparrow bandwidth h

Samples synthetic \underline{w} and y



Observe real data



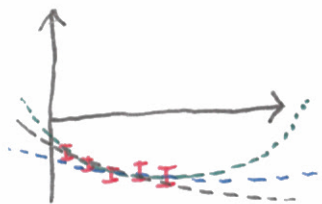
Model posterior:

$$p(h | y, X) \propto \underbrace{p(y | X, h)} p(h)$$

Marginal likelihood

Picking $h=0.05$ vs $h=1$ is just a Gaussian classifier!

Posterior for $h=1$

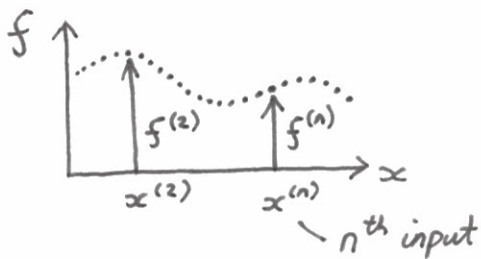


$$p(\underline{w} | y, X, h=1) = N(\underline{w}; \underline{w}_N, V_N)$$

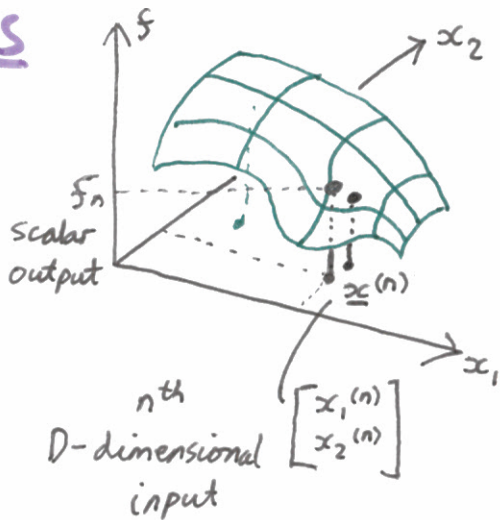
samples of $f(x; \underline{w})$

(Variance of f at \underline{x} : $\underline{x}^T V_N \underline{x}$)

Gaussian Processes



$$\underline{f} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_N \end{bmatrix} \quad \begin{array}{l} N \text{ outputs for} \\ \{ \underline{x}^{(n)} \} \text{ stored} \\ \text{in } X. \end{array}$$



*covariance /
kernel function*

If function $f \sim \text{GP}(m, k)$ \Rightarrow $P(\underline{f}) = \mathcal{N}(\underline{f}; \underline{m}, K)$

mean function \uparrow

$$m_i = m(\underline{x}^{(i)}) \text{ usually } 0$$

$$K_{ij} = k(\underline{x}^{(i)}, \underline{x}^{(j)})$$

\hookrightarrow Mercer kernels

Example:

$$k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \exp(-\|\underline{x}^{(i)} - \underline{x}^{(j)}\|^2)$$

We need k to always give positive definite k
semi-

Things we can do with Gaussians

For a joint Gaussian:

$$p(\underline{f}, \underline{g}) = N\left(\begin{bmatrix} \underline{f} \\ \underline{g} \end{bmatrix}; \begin{bmatrix} \underline{a} \\ \underline{b} \end{bmatrix}, \begin{bmatrix} A & C \\ C^T & B \end{bmatrix}\right)$$

Marginals

$$\begin{aligned} p(\underline{f}) &= \int p(\underline{f}, \underline{g}) d\underline{g} \\ &= N(\underline{f}; \underline{a}, A) \end{aligned}$$

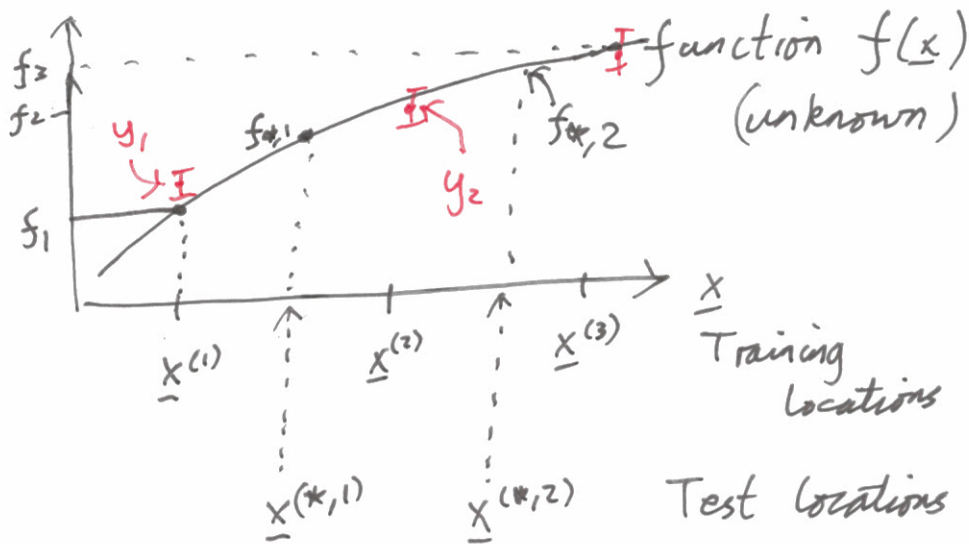
Conditionals

$$p(\underline{f} | \underline{g}) = N(\underline{f}; \underline{a} + CB^{-1}(\underline{g} - \underline{b}), A - CB^{-1}C^T)$$

also

$$p(\underline{g} | \underline{f}) = N(\underline{g}; \underline{b} + C^T A^{-1}(\underline{f} - \underline{a}), B - C^T A^{-1}C)$$

Apply to Regression



$$P(\underline{f}, \underline{f}_*) = N \left(\begin{bmatrix} \underline{f} \\ \underline{f}_* \end{bmatrix}, \begin{bmatrix} \underline{0} \\ \underline{0} \end{bmatrix}, \begin{bmatrix} k(x, x) & k(x, x_*) \\ k(x_*, x) & k(x_*, x_*) \end{bmatrix} \right)$$

Notation

$$K(X, Z)_{ij} = k(x^{(i)}, z^{(j)})$$

\uparrow \uparrow \uparrow
 $N \times D$ $M \times D$
 $N \times M$

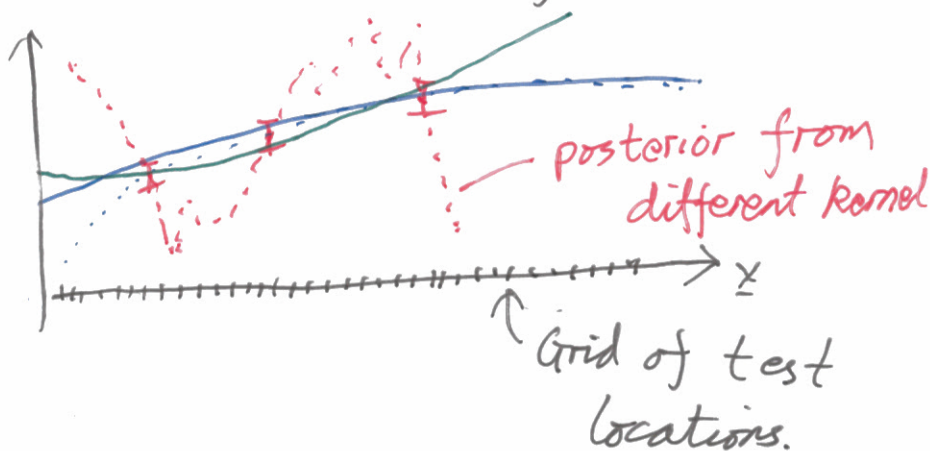
$$p(\underline{y}, \underline{f}_*)$$

$$= N \left(\begin{bmatrix} \underline{y} \\ \underline{f}_* \end{bmatrix}; \begin{bmatrix} \underline{0} \\ \underline{0} \end{bmatrix}, \begin{bmatrix} k(x, x) + \sigma_y^2 \mathbb{I} & k(x, x_*) \\ k(x_*, x) & k(x_*, x_*) \end{bmatrix} \right)$$

Inference

$$p(\underline{f}_* | \underline{y}) = N(\underline{f}_* | \underline{\mu}, \underline{\Sigma})$$

From standard results
for Gaussians.



Bayesian Linear Regression is a GP

Model

$$f_i = f(\underline{x}^{(i)}) = \underline{w}^T \underline{x}^{(i)} + b$$

Prior $\underline{w} \sim \mathcal{N}(\underline{0}, \sigma_w^2 \mathbb{I}), b \sim \mathcal{N}(0, \sigma_b^2)$

$$\begin{aligned} \text{cov}(f_i, f_j) &= \mathbb{E}[f_i f_j] - \underbrace{\mathbb{E}[f_i] \mathbb{E}[f_j]}_0 \\ &= \mathbb{E}[(\underline{w}^T \underline{x}^{(i)} + b)^T (\underline{w}^T \underline{x}^{(j)} + b)] \\ &= \mathbb{E}[\underline{x}^{(i)T} \underline{w} \underline{w}^T \underline{x}^{(j)} + b^2 + \dots] \\ &= \underbrace{\underline{x}^{(i)T} \mathbb{E}[\underline{w} \underline{w}^T] \underline{x}^{(j)}}_{\sigma_w^2 \mathbb{I}} + \underbrace{\mathbb{E}[b^2]}_{\sigma_b^2} + \underbrace{\dots}_0 \\ k(\underline{x}^{(i)}, \underline{x}^{(j)}) &= \sigma_w^2 \underbrace{\underline{x}^{(i)T} \underline{x}^{(j)}}_{\phi(\underline{x}^{(i)})^T \phi(\underline{x}^{(j)})} + \sigma_b^2 \end{aligned}$$