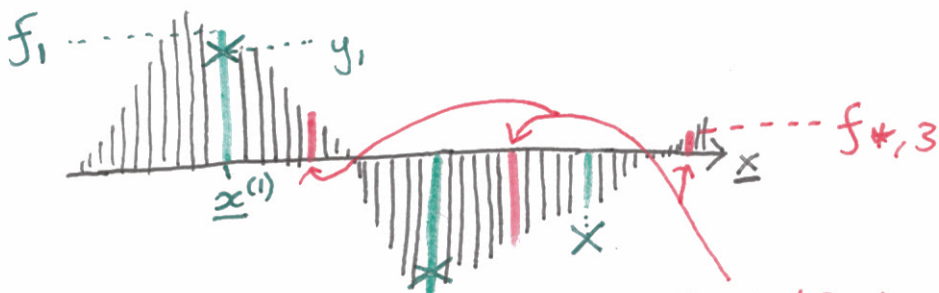


Gaussian Processes

"Function viewpoint"

Function values are parameters of model



Prior on function

$$GP(k)$$

Test / Probe locations X_*

$$\Rightarrow p(\underline{f}, \underline{f}_*) = N\left(\begin{bmatrix} \underline{f} \\ \underline{f}_* \end{bmatrix}; \underline{0}, \begin{bmatrix} k_{xx} & k_{x_*} \\ k_{x_*} & k_{**} \end{bmatrix}\right)$$

Posterior

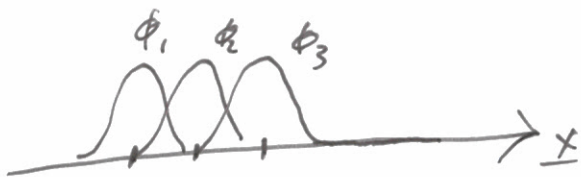
$$p(\underline{f}_* | \underline{y}) = N(\dots) \quad (\text{see notes})$$

"Weight Space view"

Prior $\underline{w} \sim N(\underline{0}, \sigma_w^2 \mathbf{I})$

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

$$\Rightarrow k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma_w^2 \underline{\phi}(\underline{x}^{(i)})^T \underline{\phi}(\underline{x}^{(j)})$$



If we put RBF everywhere

P84

GPMK

→ Analytically derive $\underline{\phi}^T \underline{\phi}$

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

"Kernel trick"

- Rewrite algorithm so it only depends on inner products

- Set $\underline{\phi}(\underline{x}^{(i)})^T \underline{\phi}(\underline{x}^{(j)}) = \underbrace{k(\underline{x}^{(i)}, \underline{x}^{(j)})}_{\text{"Mercer kernel"}}$

"Mercer kernel"

Aside Kernel Logistic Regression

(Non-exam.)

cost

Linear case, SGD, "one epoch" / $O(DN)$

$$\left[\begin{array}{l} \underline{w} \leftarrow \underline{0} \quad D \times 1 \quad \text{step size } \lambda \\ \text{for } n=1 \dots N: \\ \quad \underline{w} \leftarrow \underline{w} + u^{(n)} \underline{x}^{(n)} \end{array} \right. \begin{array}{l} \text{gradient} \\ \text{involves} \\ y^{(n)} \& \text{ confiden.} \\ - \lambda \underline{w} \end{array}$$

\underline{w} is in span of $\{\underline{x}^{(n)}\}$

λ regularizes
 \times step size

$$\underline{w} = \sum_n a_n \underline{x}^{(n)} = X^T \underline{a}$$

$$\left[\begin{array}{l} \underline{a} \leftarrow \underline{0} \quad N \times 1 \\ \text{for } n=1 \dots N: \\ \quad \underline{a} \leftarrow (1-\lambda) \underline{a} \\ \quad a_n \leftarrow a_n + u^{(n)} \end{array} \right. \begin{array}{l} \text{orig} \\ \text{dim.} \\ O(N^2 D) \end{array}$$

To get $u^{(n)}$ or to predict

$$\begin{aligned} \text{need } \underline{w}^T \underline{x}^{(*)} &= \sum_n a_n \underline{x}^{(n)T} \underline{x}^{(*)} \\ &= \sum_n a_n k(\underline{x}^{(n)}, \underline{x}^{(*)}) \end{aligned}$$

For 1 prediction:

$$p(f_* | \underline{y}) = N(f_*; m, s^2)$$

Can show: $N \times N$ Noise variance

$$m = \underline{k}^{*\top} \left(\underline{k}(\underline{x}, \underline{x}) + \sigma_y^2 \mathbb{I} \right)^{-1} \underline{y}$$

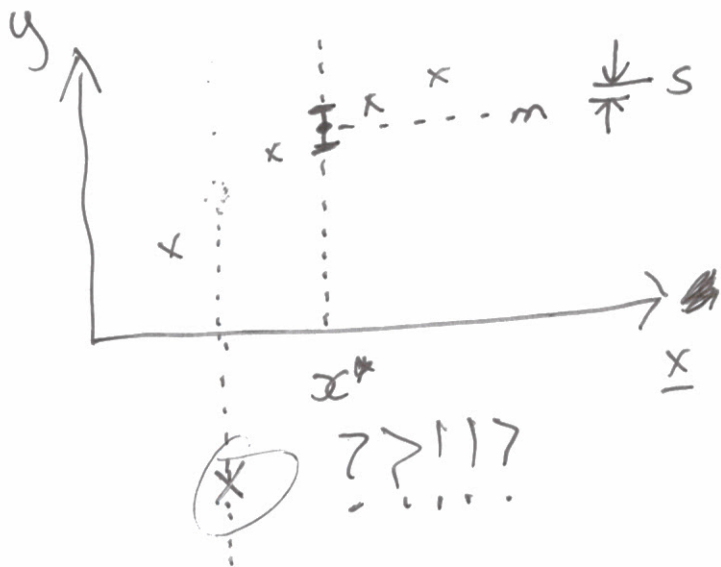
\uparrow $(\underline{k}^*)_i = k(\underline{x}^*, \underline{x}^{(i)})$ training labels

$$s^2 = k(\underline{x}^{(*)}, \underline{x}^{(*)}) - \underline{k}^{*\top} \left(\underline{k}(\underline{x}, \underline{x}) + \sigma_y^2 \mathbb{I} \right)^{-1} \underline{k}^{(*)}$$

positive definite

Has no \underline{y} dependence.

positive



Can be more uncertain at x^*
in response to surprises?

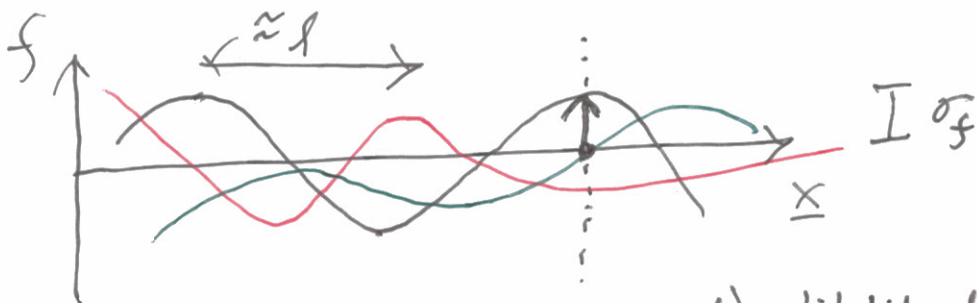
→ Change the GP kernel function.

⇒ Learn the kernel... how?

Learning the kernel

$$k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_d (x_d^{(i)} - x_d^{(j)})^2 / \ell_d^2\right)$$

Also learn noise level σ_y^2



Pick parameters by (marginal) likelihood:

$$p(y | X, \theta = \{\sigma_y^2, \sigma_f^2, \{\ell_d\}\})$$

$$= N(y; \underline{0}, k(x, x) + \sigma_y^2 \mathbf{I})$$