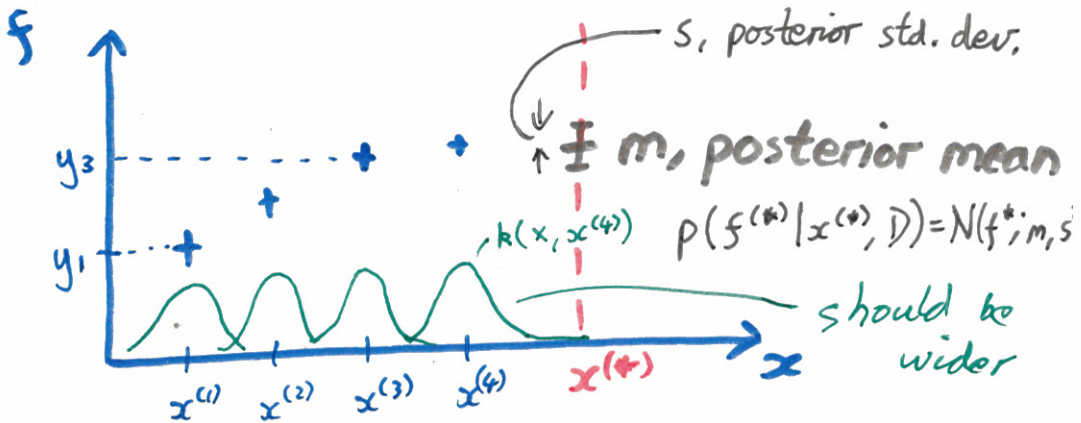


GAUSSIAN PROCESSES



$N=4$ training points

$N \times N$ matrix K , $k_{ij} = k(x^{(i)}, x^{(j)})$

Can show: $(\underline{k}^{(*)})_i = k(x^{(*)}, x^{(i)})$

$$m = \underline{k}^{(*)T} (K + \sigma_n^2 I)^{-1} \underline{y}$$

$$s^2 = \underbrace{k(x^{(*)}, x^{(*)})}_{\text{Usually } \sigma_f^2} - \underbrace{\underline{k}^{(*)T} (K + \sigma_n^2 I)^{-1} \underline{k}^{(*)}}_{M^{-1}}$$

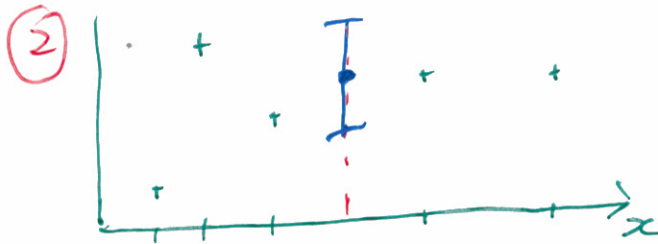
$$k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma_f^2 e^{-\frac{1}{2} \|\underline{x}^{(i)} - \underline{x}^{(j)}\|} \quad \text{Has no dependence on } \underline{y}$$



Long lengthscale

Short lengthscale.

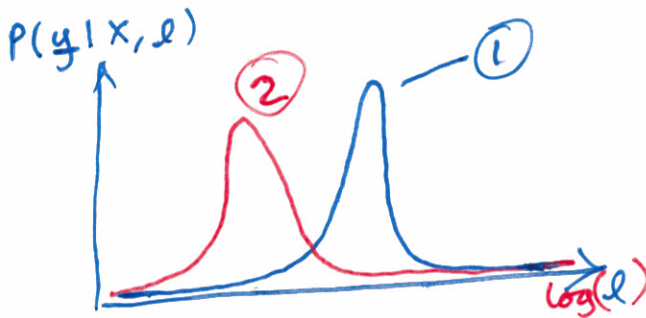
Should be more uncertain if y are surprising.



We need learn params of kernel

$$k(x^{(i)}, x^{(j)}) = e^{-\frac{1}{2}(x^{(i)} - x^{(j)})^2 / l^2}$$

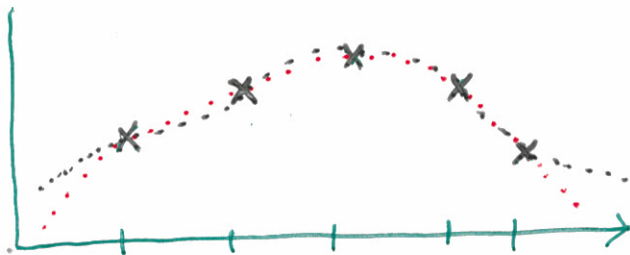
↑
lengthscale



Short l
Long l

Assuming I know small noise level.

Look at the code demo on the website



Bayesian Logistic Regression

$$P(y=1 | \underline{x}, \underline{w}) = \sigma(\underline{w}^T \underline{x}) = \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$$

Maximized likelihood:

$$\begin{aligned} P(\underline{y} | X, \underline{w}) &= \prod_n \sigma(z^{(n)} \underline{w}^T \underline{x}) && z^{(n)} \in \{\pm 1\} \\ &= P(D | \underline{w}) && z^{(n)} = 2y^{(n)} - 1 \\ &\quad \uparrow && y^{(n)} \in \{0, 1\} \\ &\quad \text{Training data.} \end{aligned}$$

Minimize a penalized neg. log likelihood

$$\underline{w}^* = \underset{\underline{w}}{\operatorname{arg\,min}} \left[-\log P(D | \underline{w}) + \lambda \underline{w}^T \underline{w} \right]$$

Fit, guess of weights.

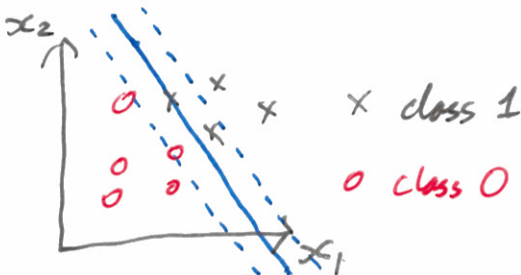
Bayes' Rule

$$p(\underline{w} | D) = \frac{p(D | \underline{w}) p(\underline{w})}{p(D)} \propto p(D | \underline{w}) p(\underline{w})$$

$$p(D) = \int p(D | \underline{w}) p(\underline{w}) d\underline{w} \quad \left. \begin{array}{l} \text{we can't} \\ \text{do this} \\ \text{integral exactly.} \end{array} \right\}$$

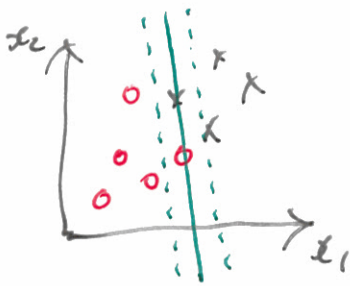
$p(D)$ or $p(D | M)$

↳ model, choice of basis f^n 's
choice of prior



Decision boundary for one setting of weights

$$\sigma(\underline{w}^T \underline{x}) = 0.27 \quad \sigma(\underline{w}^T \underline{x}) = 1/2$$

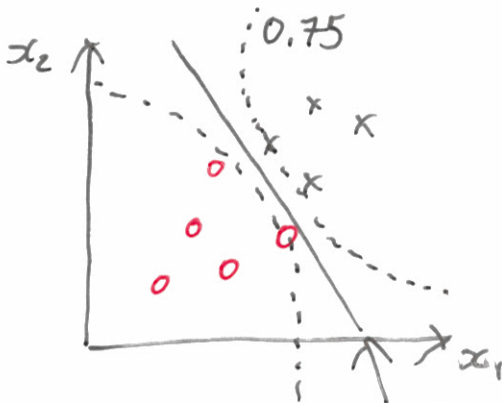


--- contour of
 $p(y=1 | \underline{x}, \underline{w})$

Bayesian Predictions

$$p(y | \underline{x}, D) = \int p(y, \underline{w} | \underline{x}, D) d\underline{w} \quad (\text{sum rule})$$

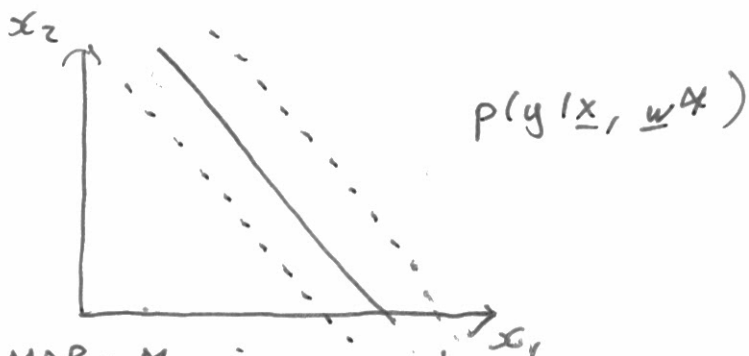
$$= \int p(y | \underline{x}, \underline{w}) p(\underline{w} | D) d\underline{w} \quad (\text{product rule})$$



Contours
 $p(y | \underline{x}, D)$

0.5, decision boundary

Compare to MAP predictions



MAP = Maximum a posteriori:

$$\underline{w}^* = \underset{\underline{w}}{\operatorname{argmax}} \log p(\underline{w} | D)$$

$$= \underset{\underline{w}}{\operatorname{argmax}} \left[\log p(D | \underline{w}) + \underbrace{\log p(\underline{w})}_{-\frac{1}{2\sigma_w^2} \underline{w}^T \underline{w} + \text{const}} \right]$$

If $p(\underline{w}) = N(\underline{w}; 0, \sigma_w^2 I)$

MAP is not Bayesian,

can be seen as a crude approximation