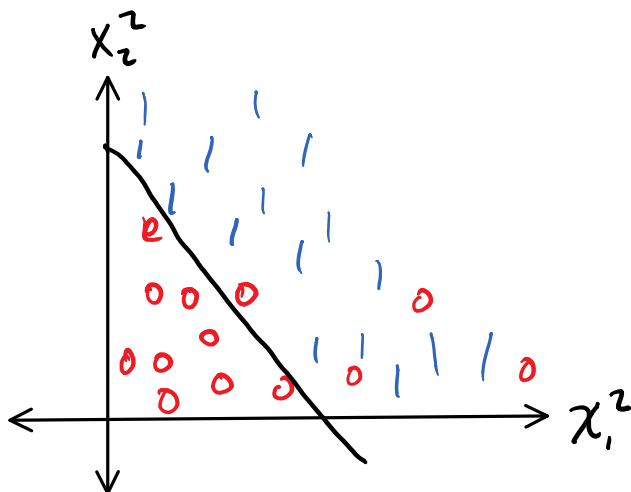


Linear decision
boundary

Non-linear decision
boundary

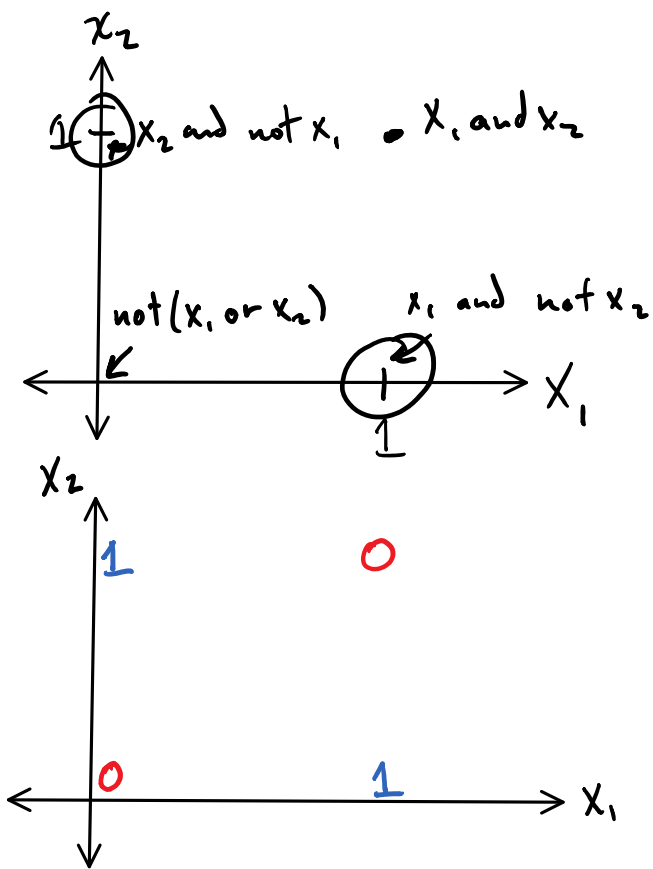
define higher dim embedding $\Phi: \mathbb{R}^P \rightarrow \mathbb{R}^D$
 $\Phi(x) \in \mathbb{R}^D$

ex $\Phi(x_1, x_2) = (1, x_1, x_2, x_1^2, x_2^2)$



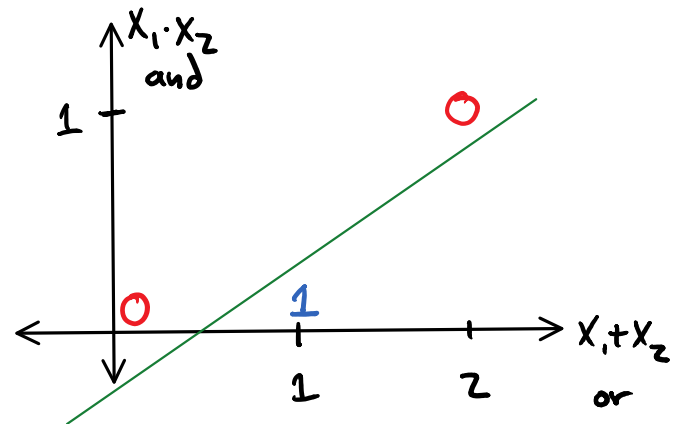
Φ make linear methods
into non-linear method

ex Logic: x_1, \dots, x_p are propositions encoded as $\{0,1\}$.



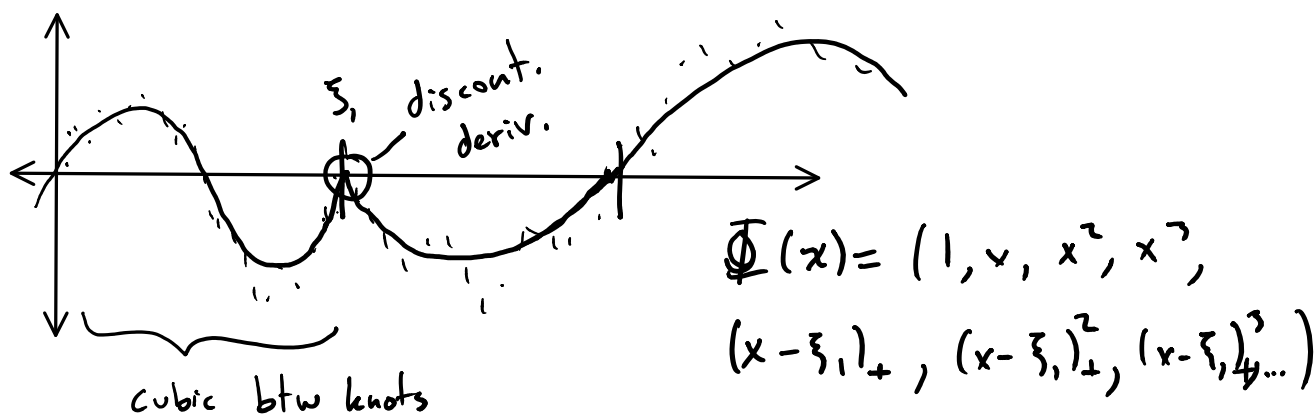
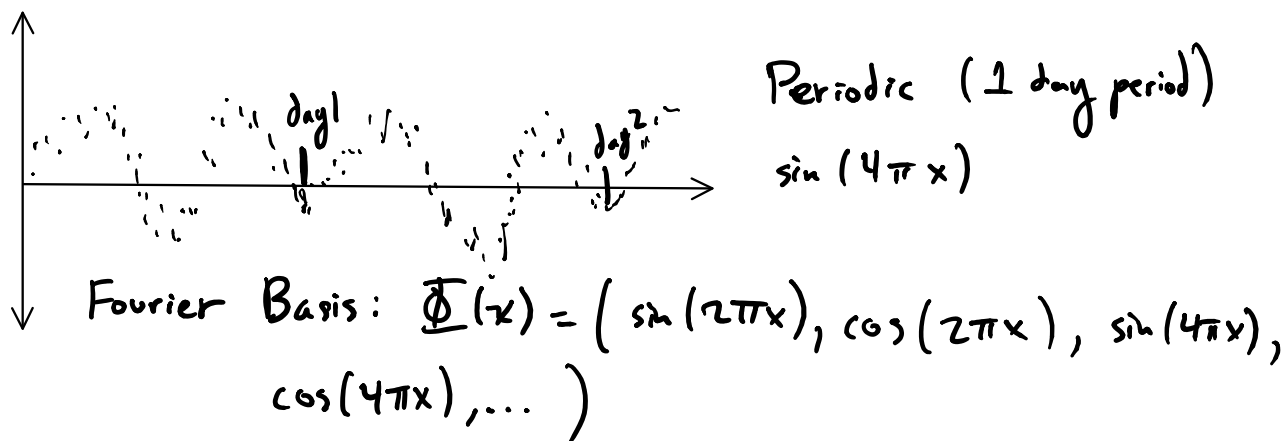
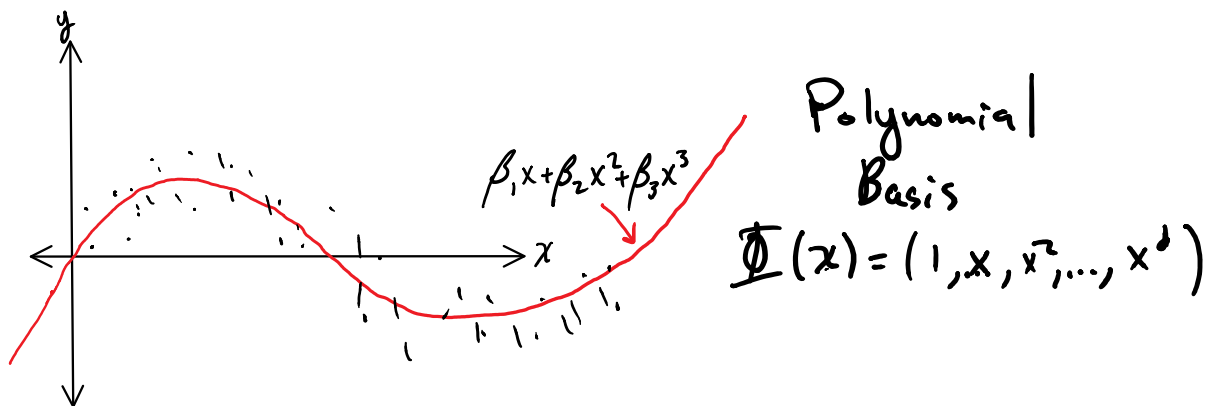
$$\frac{x_1 \text{ xor } x_2}{(x_1 \text{ and not } x_2) \text{ or } (x_2 \text{ and not } x_1)}$$

$$\Phi(x_1, x_2) = (x_1, x_2, x_1 \oplus x_2)$$



Basis Expansion

Monday, May 8, 2017 8:42 PM

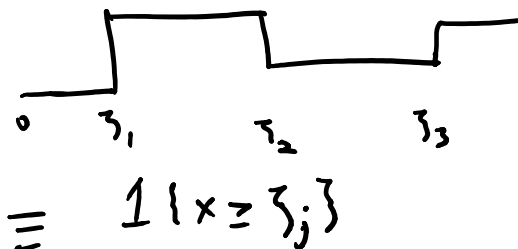


Derivatives & Constraints

0th order: $\phi_1(x) = 1 \{0 \leq x < \xi_1\}$

$\phi_2(x) = 1 \{\xi_1 \leq x < \xi_2\}$

$\phi_3(x) = 1 \{\xi_2 \leq x < \xi_3\}$



k^{th} order poly: $1, (x - \xi_1)_+, (x - \xi_1)_+^2, \dots, (x - \xi_1)_+^k$
 $k-1$ cont deriv: then remove

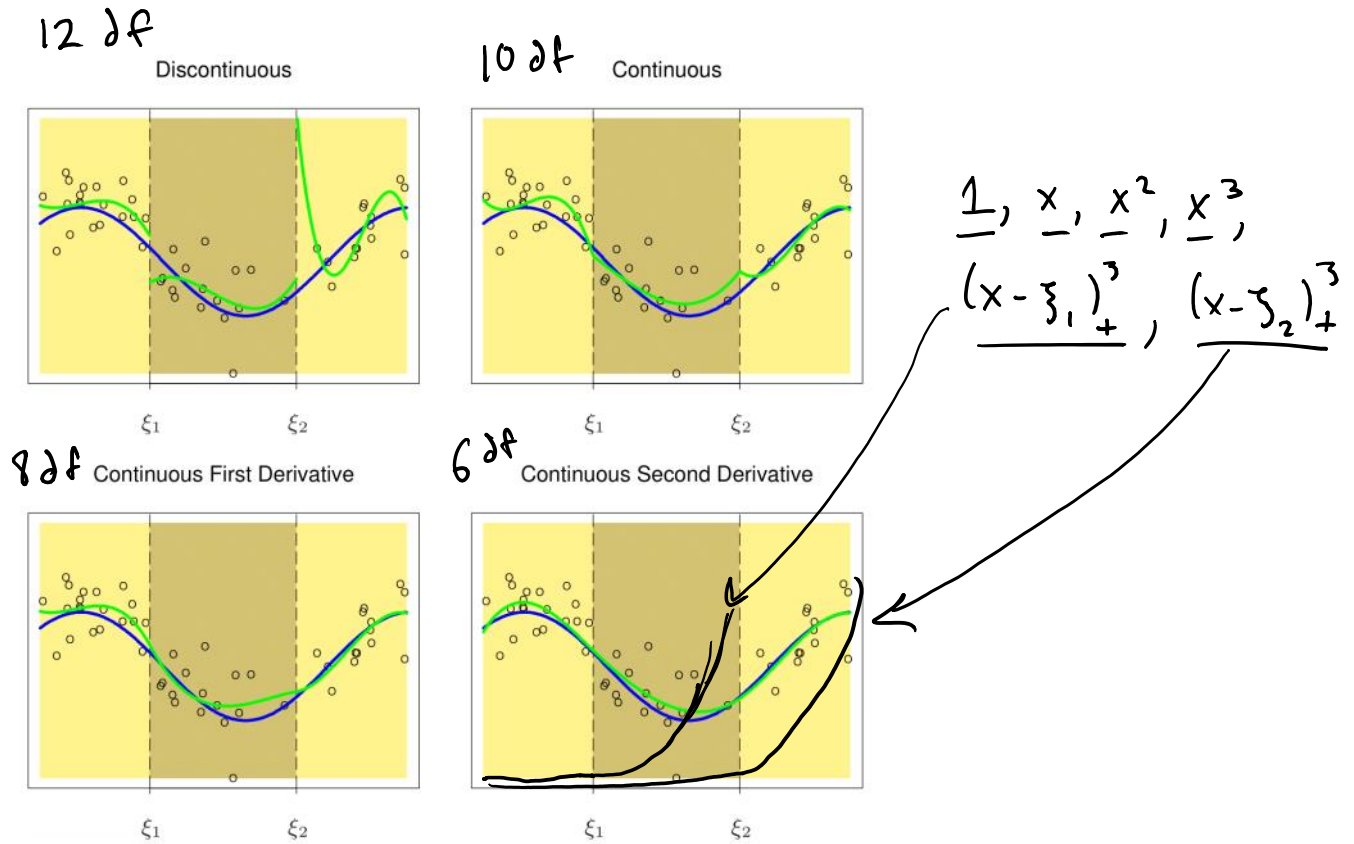


FIGURE 5.2. A series of piecewise-cubic polynomials, with increasing orders of continuity.

ESL 5.2

Kernel Trick

let $z_{i,l} = \phi_l(x_i) \quad i=1, \dots, n \quad l=1, \dots, D$

SVM for $y_i \in \{-1, 1\}$

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n (1 - y_i z_i^T \beta)_+ + \lambda \|\beta\|_2^2$$

(I) (II)

$z\beta \leftrightarrow \text{loss}$

claim $\hat{\beta}$ solves SVM can be written as $z^T \alpha$
 $\alpha \in \mathbb{R}^D$ i.e. $\hat{\beta}_j = z_j^T \alpha = \sum_i \alpha_i z_{ij}$

proof $\beta = \sum_i \alpha_i z_i + \beta^\perp \quad z_i^T \beta^\perp = 0$ (I)

$$= z_i^T z^T \alpha$$

β^\perp does not impact R_n

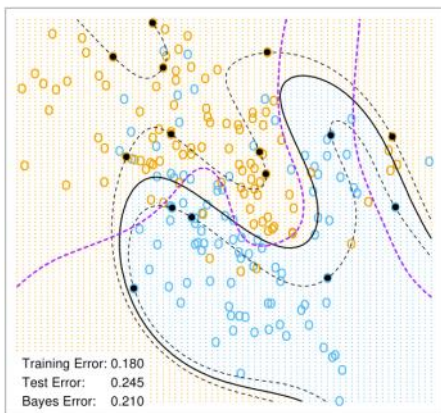
$$\|\beta\|_2^2 = \|z^T \alpha + \beta^\perp\|_2^2 = \|z^T \alpha\|_2^2 + 2 \underbrace{\beta^\perp^T z^T \alpha}_{(z\beta^\perp)^T \alpha = 0} + \|\beta^\perp\|_2^2$$

$$= \|z^T \alpha\|_2^2 + \|\beta^\perp\|_2^2 \quad \text{(II)}$$

SVM: $\min_{\alpha \in \mathbb{R}^D, \beta^\perp \in \mathbb{R}^D} \frac{1}{n} \sum_i (1 - y_i z_i^T z^T \alpha)_+ + \lambda (\|z^T \alpha\|_2^2 + \|\beta^\perp\|_2^2)$
 s.t. $\beta^\perp^T z_i = 0 \quad \forall i$

min'ed when $\beta^\perp = 0$

SVM - Degree-4 Polynomial in Feature Space



SVM - Radial Kernel in Feature Space



Ridge Regression

$$\min_{\beta} \sum_i (y_i - z_i^T \beta)^2 + \lambda \|\beta\|_2^2$$

same story!

General

$$\min_{\beta \in \mathbb{R}^D} R_n(y, z\beta) + \lambda \|\beta\|_2^2$$

$$\min_{\alpha \in \mathbb{R}^n} R_n(y, z z^T \alpha) + \lambda \underbrace{\|z^T \alpha\|_2^2}_{\alpha^T z z^T \alpha}$$

define $K = z z^T \quad (K_{ij} = z_i^T z_j = \Phi(x_i)^T \Phi(x_j))$

$$\min_{\alpha} R_n(y, K\alpha) + \lambda \alpha^T K \alpha$$

Method 1 define transformation

Φ compute either $z = \Phi(x)$
 or compute $K = \Phi(x)^T \Phi(x)$
 and solve SVM

Method 2 define kernel function

$$k(x_i, x_j) := \Phi(x_i)^T \Phi(x_j)$$

then compute K and solve SVM

Do we need Φ to use a kernel k ?

ex $k(x, x') = e^{-\frac{\|x - x'\|_2^2}{2\sigma^2}}$
 σ bandwidth parameter.

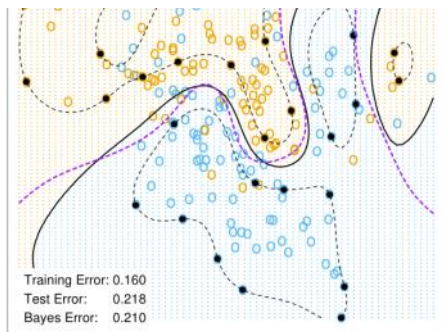


FIGURE 12.3. Two nonlinear SVMs for the mixture data. The upper plot uses a 4th degree polynomial kernel, the lower a radial basis kernel (with $\gamma = 1$). In each case C was tuned to approximately achieve the best test error performance, and $C = 1$ worked well in both cases. The radial basis kernel performs the best (close to Bayes optimal), as might be expected given the data arise from mixtures of Gaussians. The broken purple curve in the background is the Bayes decision boundary.

ESL 12.3

def Mercer kernel is a function

$$k: \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}_+ \text{ that is PSD}$$

(for any $\{x_i\} \subseteq \mathbb{R}^d$ $(k(x_i, x_j))_{ij}$ is PSD)

ex d^{th} degree poly. $k(x, x') = (1 + x^T x')^d$

$$\begin{aligned} d=2: (1 + x_1 x'_1 + x_2 x'_2)^2 &= 1 + 2x_1 x'_1 + 2x_2 x'_2 \\ &+ x_1^2 x'^2_1 + x_2^2 x'^2_2 + 2x_1 x'_1 x_2 x'_2 \end{aligned}$$

$$= (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2)^T \cdot (\dots x' \dots)$$

$$\hookrightarrow \Phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2, x_1^2, x_2^2)$$

thm Every Mercer kernel has a Hidi embedding

$$\Phi \text{ s.t. } k(x, x') = \Phi(x)^T \Phi(x')$$

(Φ may be infinite dim)

ex RBF in 1d

$$\Phi(x) = e^{-x^2/2\sigma^2} [1, \sqrt{\frac{1}{2!\sigma^2}} x, \sqrt{\frac{1}{4!\sigma^4}} x^2, \dots]$$

Predict new x^*

$$\Phi(x^*)^T \hat{\beta} = \Phi(x^*)^T \underset{\uparrow \text{training}}{Z^T} \hat{\alpha} = \sum_i \hat{\alpha}_i (\Phi(x^*)^T \Phi(x_i))$$

$$= \sum_i \hat{\alpha}_i k(x^*, x_i)$$

- now predict scales w/ n!