

lecture 12

- C-SVM \rightarrow selecting C
- Non-linear SVM
 - Kernel trick
 - Nonlinear SVM prediction
 - RBF - effect of σ

Project
submit clustering
results alg1.txt
~~alg2.txt~~
5/23 \leftarrow = 10% - 15%
Report : exam week
 $\approx 90\% - 85\%$

Lecture V: Support Vector Machines and Kernel Machines

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

April, 2023

Linear SVM's

The margin and the expected classification error

Maximum Margin Linear classifiers

Linear classifiers for non-linearly separable data

C-SVM

Non linear SVM

The “kernel trick”

Kernels

Prediction with SVM

Extensions

L_1 SVM

Multi-class and One class SVM

SV Regression

Reading AoNPS Ch.: Ch. 12.1–3, HTF Ch.: Ch 14 (14.1,14.2–14.2.4 kernels, 14.4 and equations (14.28,14.29) kernel trick, 14.5.1.–3 Support Vector Machines) 7.1–7.4, 7.7

Additional Reading: C. Burges - “A tutorial on SVM for pattern recognition”

These notes: Appendices (convex optimization) are optional.

Non-linearly separable problems and their duals

The C-SVM

$$\begin{aligned} & \underset{w, b, \xi}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y^i (w^T x^i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \quad \text{slack} \end{aligned} \tag{20}$$

regularization

In the above, ξ_i are the **slack variables**. Dual³:

$$\begin{aligned} & \underset{\alpha}{\text{maximize}} \quad \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^i y^j x^{iT} x^j \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0 \text{ for all } i \\ & \sum_i \alpha_i y^i = 0 \end{aligned} \tag{21}$$

$\alpha^T G \alpha \approx G_{ij}$

⇒ two types of SV

- $\alpha_i < C$ data point x^i is “on the margin” $\Leftrightarrow y^i(w^T x^i + b) = 1$ (original SV)
- $\alpha_i = C$ data point x^i cannot be classified with margin 1 (**margin error**)
 $\Leftrightarrow y^i(w^T x^i + b) < 1$

³Lagrangian $L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y^i(w^T x^i + b) - 1 + \xi_i] - \sum_i \mu_i \xi_i$ with $\alpha_i \geq 0, \xi_i \geq 0, \mu_i \geq 0$

The margin and the expected classification error

Theorem Let $\mathcal{F} = \{\text{sgn}(w^T x), \|w\| \leq \Lambda, \|x\| \leq R\}$ and let $\rho > 0$ be any “margin”. Then for any $f \in \mathcal{F}$, w.p $1 - \delta$ over training sets

$$\Pr[\text{err}] \equiv L_{01}(f) \leq \hat{L}_\rho + \sqrt{\frac{c}{n} \left(\frac{R^2 \Lambda^2}{\rho^2} \ln n^2 + \ln \frac{1}{\delta} \right)} \quad \begin{matrix} \leftarrow \text{find } \rho \\ \text{that minimizes} \end{matrix} \quad (5)$$

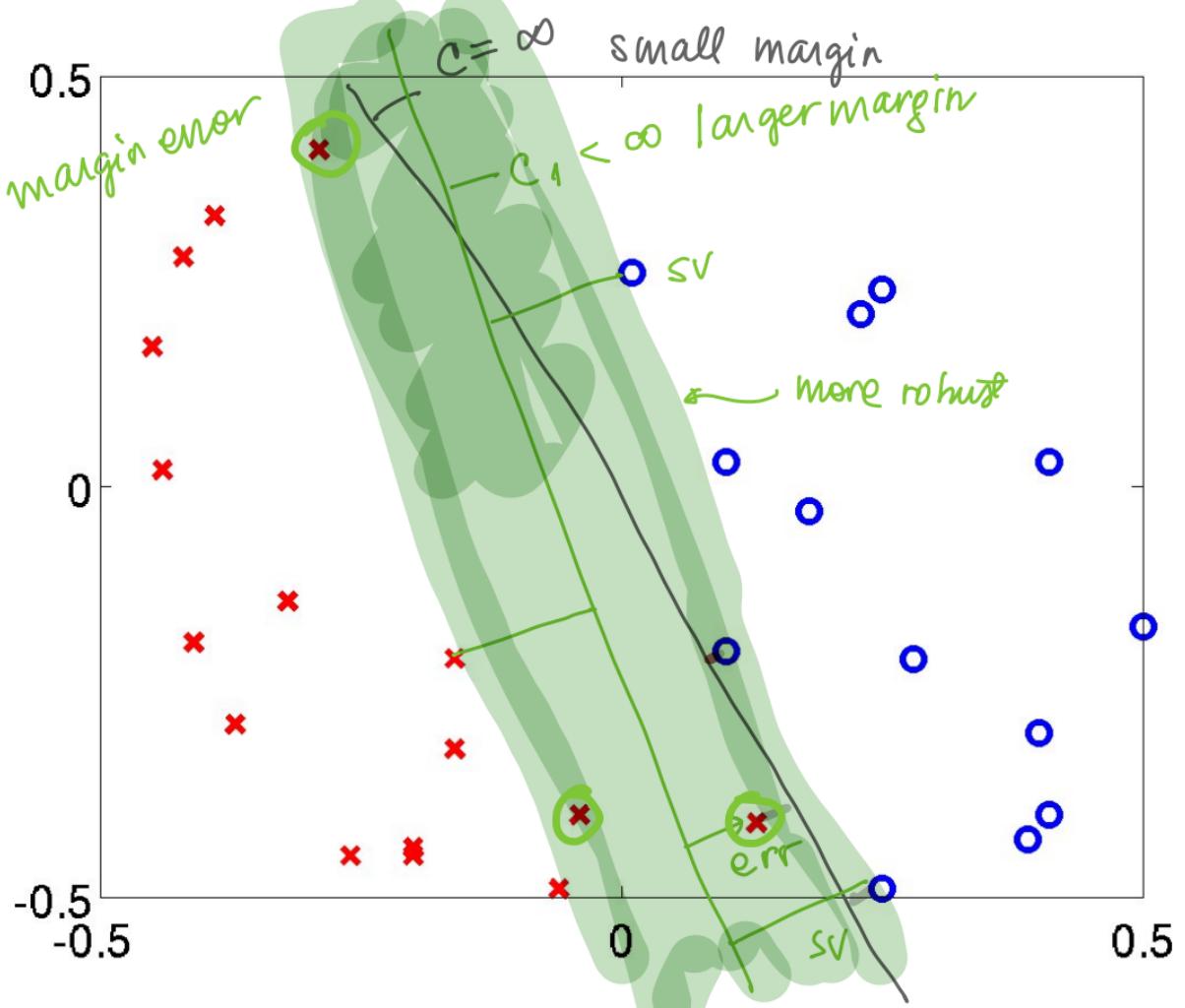
where c is a universal constant and \hat{L}_ρ is the fraction of the training examples for which

$$y^i w^T x_i < \rho \quad (6)$$

- ▶ a data point i that satisfies (6) for some ρ is called a **margin error**
- ▶ For $\rho = 0$ the margin error rate \hat{L}_ρ is equal to \hat{L}_{01}

Choosing C

- $C \uparrow \Rightarrow \sum_{i=1}^n \xi_i \rightarrow$ penalized more \Rightarrow fewer & margin errors smaller
 $\Rightarrow \|w\| \uparrow \Leftrightarrow \text{margin } \Downarrow$
- try $C \Rightarrow$ for all $\hat{\xi}_p = \frac{\# \text{margin errors } (\hat{\xi})}{n}$ $\rightarrow \min_p \text{r.h.s (5)} \equiv \text{Bound}(C)$
(for $P[\text{err}]$)
compute $\text{SVM}(C)$
 $\min C \text{ Bound}(C)$
- $\hat{v} = \frac{1}{C}$ (heuristic)
- Cross Validation



The ν -SVM

regularization

$$\underset{w, b, \xi, \rho}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 - \nu\rho + \frac{1}{n} \sum_i \xi_i \quad (22)$$

$$\text{s.t.} \quad y^i(w^T x^i + b) \geq \underline{\rho} - \underline{\xi}_i \quad (23)$$

$$\xi_i \geq 0 \quad (24)$$

$$\rho \geq 0 \quad (25)$$

where $\nu \in [0, 1]$ is a parameter.Dual⁴:

n small
n large

$$\underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_i \alpha_i \alpha_j y^i y^j x^{iT} x^j \quad (26)$$

$$\text{s.t.} \quad \frac{1}{n} \geq \alpha_i \geq 0 \text{ for all } i \quad (27)$$

$$\sum_i \alpha_i y^i = 0 \quad (28)$$

$$\sum_i \alpha_i \geq \nu \quad (29)$$

Properties If $\rho > 0$ then:

- ν is an upper bound on #margin errors/n (if $\sum_i \alpha_i = \nu$)
- ν is a lower bound on #(original support vectors + margin errors)/n
- ν -SVM leads to the same w, b as C-SVM with $C = 1/\nu$

⁴Lagrangian $L(w, b, \xi, \rho, \alpha, \mu, \delta) = \frac{1}{2} \|w\|^2 - \nu\rho + \frac{1}{n} \sum_i \xi_i - \sum_i \alpha_i [y^i(w^T x^i + b) - \rho + \xi_i] - \sum_i \mu_i \xi_i - \delta\rho$
with $\alpha_i \geq 0, \delta \geq 0, \mu_i \geq 0$

A simple error bound

$$L_{01}(f_n) \leq E \left[\frac{\#\text{support vectors of } f_{n+1}}{n+1} \right] \quad (30)$$

where f_n denotes the SVM trained on a sample of size n .

Exercise Use the Homework 6 to prove this result.

Non-linear SVM

How to use linear classifier on data that is not linearly separable?

An old trick

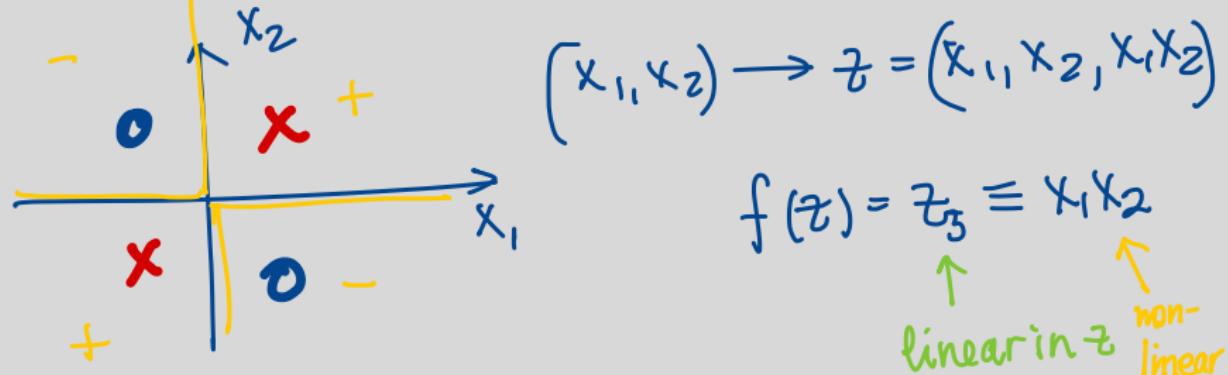
- Map the data $x^{1:n}$ to a higher dimensional space

$$x \rightarrow z = \phi(x) \in \mathcal{H}, \text{ with } \dim \mathcal{H} >> n.$$

- Construct a linear classifier $w^T z + b$ for the data in \mathcal{H}

In other words, we are implementing the non-linear classifier

$$f(x) = w^T \phi(x) + b = w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_m \phi_m(x) + b \quad (31)$$



Example

- ▶ Data $\{(x, y)\}$ below are not linearly separable

x	y	z
-1	-1	1
-1	1	-1
1	-1	-1
1	1	1

- ▶ We map them to 3 dimensions by

$$z = \phi(x) = [x_1 \ x_2 \ x_1 x_2].$$

- ▶ Now the classes can be separated by the hyperplane $z_3 = 0$ (which happens to be the maximum margin hyperplane). Hence,
 - ▶ $w = [0 \ 0 \ 1]$ (a vector in \mathcal{H})
 - ▶ $b = 0$
 - ▶ and the classification rule is $f(\phi(x)) = w^T \phi(x) + b$.
- ▶ If we write f as a function of the original x we get

$$f(x) = x_1 x_2$$

a quadratic classifier.

Non-linear SV problem

- ▶ Primal problem minimize $\frac{1}{2}||w||^2$ s.t $y^i(w^T \phi(x^i) + b) - 1 \geq 0$ for all i .
 - ▶ Dual problem

$$\max_{\alpha_{1:n}} \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j \underbrace{y^i y_j \phi(x^i)^T \phi(x_j)}_{\tilde{G}_{ij}} \text{ s.t. } \alpha_i \geq 0 \text{ for all } i \text{ and } \sum_i y^i \alpha_i = 0 \quad (32)$$

$$G_{ij} = \phi(x^i)^T \phi(x^j) \quad \text{and} \quad \bar{G} = \text{diag}\{\gamma^{1:n}\}^T G \text{diag}\{\gamma^{1:n}\}. \quad (33)$$

- ▶ \bar{G}_{ij} has been redefined in terms of ϕ
 - ▶ Dual problem

$$\max_{\alpha} \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T \bar{G} \alpha \quad \text{s.t. } \alpha_i \geq 0, \ y^T \alpha = 0 \quad (34)$$

- Same as (19)!

$$x \in \mathbb{R}^d$$

want to add all interactions

$$d + \frac{d(d+1)}{2}$$

d large $\Rightarrow \varphi(x)$ intractable !!

Example

$$(x_1, x_2) \mapsto (\cancel{x_1}, \cancel{\sqrt{2}x_2}, x_1^2, x_2^2, \cancel{\sqrt{2}x_1x_2}, 1) = 3$$

feature map

$$\begin{aligned} z^T z' &= \cancel{2x_1x'_1} + \cancel{2x_2x'_2} + \cancel{2x_1x_2x'_1x'_2} + x_1^2x'^2_1 + x_2^2x'^2_2 + 1 \quad \sim d^2 \\ &= (1 + x_1x'_1 + x_2x'_2)^2 \quad \sim d \end{aligned}$$

kernel trick

The “Kernel Trick”

Idea The result (34) is the celebrated **kernel trick** of the SVM literature. We can make the following remarks.

↗ faster than $(\dim \phi)$

1. The ϕ vectors enter the SVM optimization problem only through the Gram matrix, thus only as the scalar products $\langle \phi(x^i)^T \phi(x_j) \rangle$. We denote by $K(x, x')$ the function

$$K(x, x') = K(x', x) = \phi(x)^T \phi(x') \quad (35)$$

K is called the **kernel** function. If K can be computed efficiently, then the Gram matrix G can also be computed efficiently. This is exactly what one does in practice: we choose ϕ implicitly by choosing a kernel K . Hereby we also ensure that K can be computed efficiently.

2. Once G is obtained, the SVM optimization is independent of the dimension of x and of the dimension of $z = \phi(x)$. The complexity of the SVM optimization depends only on n the number of examples. This means that we can choose a very high dimensional ϕ without any penalty on the optimization cost.
3. Classifying a new point x . As we know, the SVM classification rule is

$$f(x) = w^T \phi(x) + b = \sum_{i=1}^n \alpha_i y^i \phi(x^i)^T \phi(x) = \sum_{i=1}^n \alpha_i y^i K(x^i, x) \quad (36)$$

Hence, the classification rule is expressed in terms of the support vectors and the kernel only. No operations other than scalar product are performed in the high dimensional space H .

Kernels

The previous section shows why SVMs are often called **kernel machines**. If we choose a kernel, we have all the benefits of a mapping in high dimensions, without ever carrying on any operations in that high dimensional space. The most usual kernel functions are

$$K(x, x') = (1 + x^T x')^p \quad \text{the polynomial kernel of degree } p$$

$$K(x, x') = \tanh(\sigma x^T x' - \beta) \quad \text{the "neural network" kernel}$$

$$K(x, x') = e^{-\frac{\|x-x'\|^2}{\sigma^2}} \quad \begin{aligned} &\text{the Gaussian or \textbf{radial basis function} (RBF) kernel} \\ &\text{it's } \underline{\phi} \text{ is } \infty\text{-dimensional} \end{aligned}$$

$$\uparrow \quad \varphi = ?$$

(Taylor expansion)

When is $K(\cdot, \cdot)$ a kernel?

when $\exists \varphi$ such that

$$\varphi(x)^T \varphi(x') = K(x, x')$$

Prediction with SVM

► Estimating b

- For any i support vector, $w^T x^i + b = y^i$ because the classification is tight
- Alternatively, if there are slack variables, $w^T x^i + b = y^i(1 - \xi_i)$
- Hence, $b = y^i(1 - \xi_i) - w^T x^i$
- For non-linear SVM, where w is not known explicitly, $w = \sum_j \alpha_j y^j \phi(x_j)$. Hence,
 $b = y^i(1 - \xi_i) - \sum_{j=1}^n \alpha_j y^j K(x^i, x^j)$ for any i support vector

► Given new x

$$\hat{y}(x) = \text{sgn}(w^T x + b) = \text{sgn} \left(\sum_{i=1}^n \alpha_i y^i K(x^i, x) + b \right). \quad (38)$$

Kernel SVM Summary

- Data $\mathcal{D} = \{(x^{1:n}, y^{1:n})\}$
- Choose kernel $K \Rightarrow G = [K(x^i, x^j)]_{i,j=1:n}$

$$\tilde{G} = [y^i y^j \alpha_{ij}]_{i,j=1:n}$$

- Choose C
- Solve (D)

$$\max_{\alpha} \underbrace{1^T \alpha}_{m \times n} - \frac{1}{2} \alpha^T \tilde{G} \alpha \quad \text{s.t.} \quad \begin{aligned} \sum_{i=1}^n \alpha_i y^i &= 0 \\ 0 \leq \alpha_i &\leq C \end{aligned}$$

• solution $\alpha_{1:n} \Rightarrow w = \sum_{i=1}^n \alpha_i \varphi(x^i) y^i, f(x) = \sum_{i=1}^n \alpha_i y^i \varphi(x^i)^T \varphi(x) + b$

- Classify new x

$$f_{w,b}(x) = \sum_{i=1}^n \alpha_i y^i K(x^i, x) + b$$

w never represented

\downarrow
compute b as before

$$\sum_{i=1}^n \alpha_i y^i K(x^i, x^b) + b = y^b$$

\downarrow
 $i_0 = \text{s.v.}$ solve for b

The Mercer condition

- ▶ How do we verify that a chosen K is a valid kernel, i.e that there exists a ϕ so that $K(x, x') = \phi(x)^T \phi(x')$?
- ▶ This property is ensured by a positivity condition known as the **Mercer** condition.

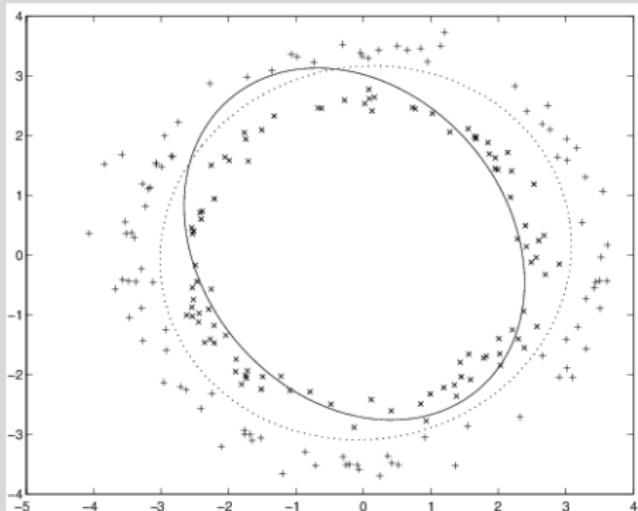
Mercer condition

Let (\mathcal{X}, μ) be a finite measure space. A symmetric function $K : \mathcal{X} \times \mathcal{X}$, can be written in the form $K(x, x') = \phi(x)^T \phi(x')$ for some $\phi : \mathcal{X} \rightarrow \mathcal{H} \subset \mathbb{R}^m$ iff

$$\int_{\mathcal{X}^2} K(x, x') g(x) g(x') d\mu(x) d\mu(x') \geq 0 \quad \text{for all } g \text{ such that } \|g(x)\|_{L_2} < \infty \quad (37)$$

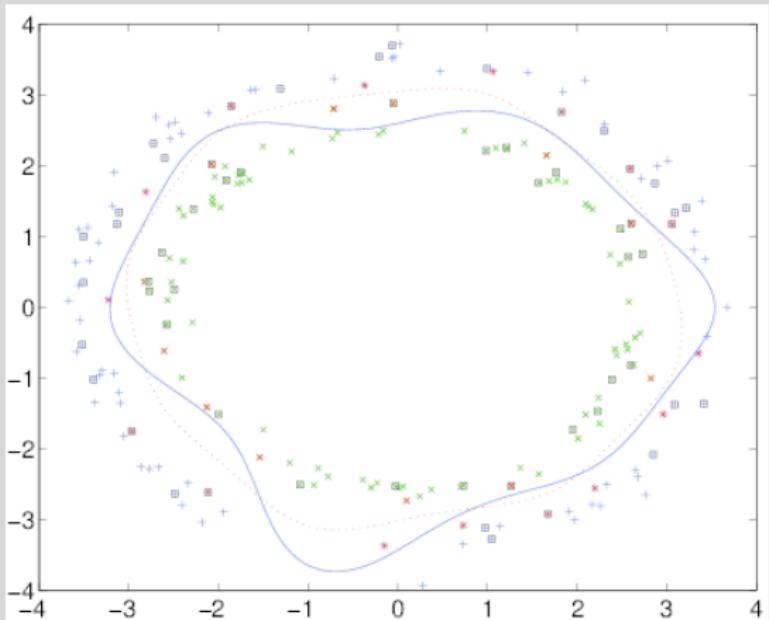
- ▶ In other words, K must be a positive semidefinite operator on L_2 .
- ▶ If K satisfies the Mercer condition, there is no guarantee that the corresponding ϕ is unique, or that it is finite-dimensional.

Quadratic kernel



- ▶ C-SVM, polynomial degree 2 kernel, $n = 200$, $C = 10000$
- ▶ The two ellipses show that a constant shift to the data ($x^i \leftarrow x^i + v$, $v \in \mathbb{R}^n$) can affect non-linear kernel classifiers.

RBF kernel and Support Vectors

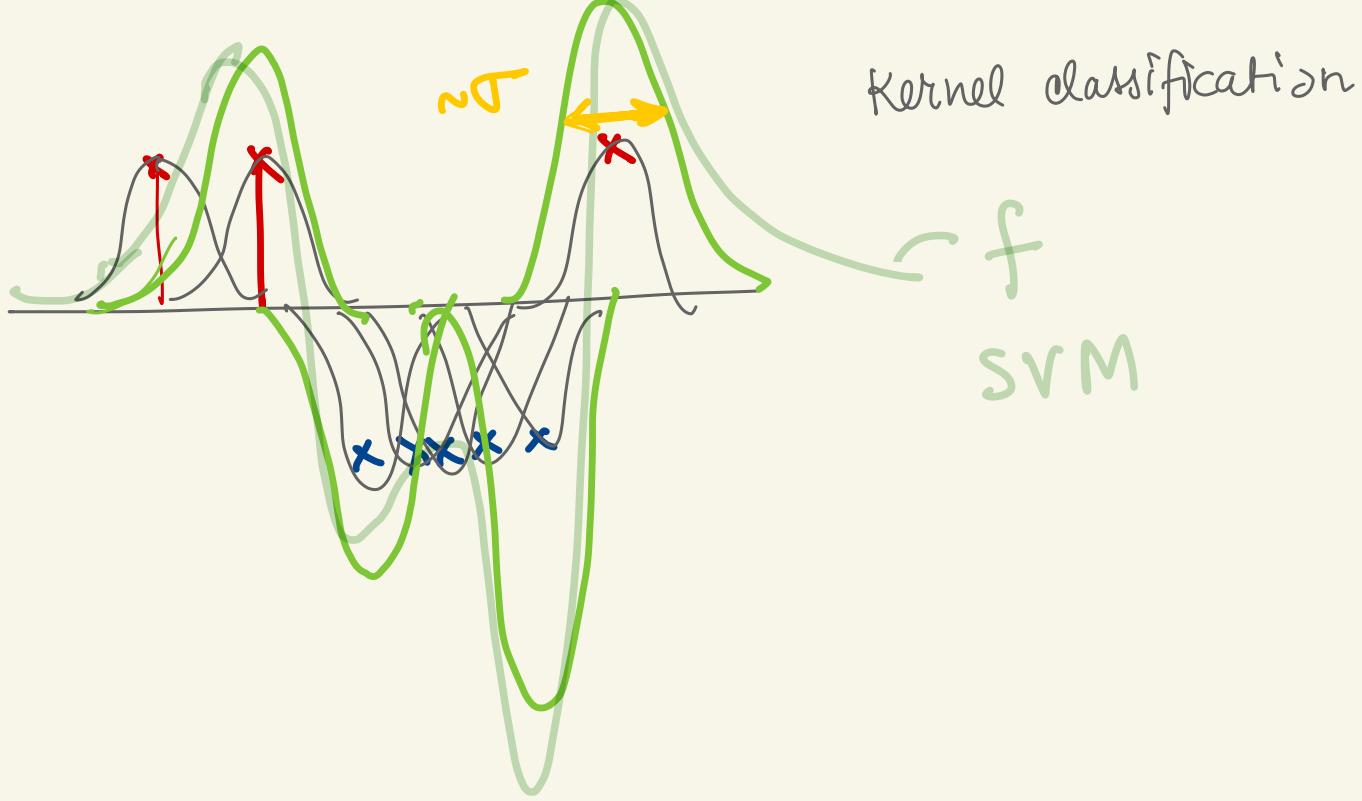


RBF Kernel

$$k(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$$

$$f(x) = \sum \alpha_i k(x^i, x) y_i + b$$

↑
RBF



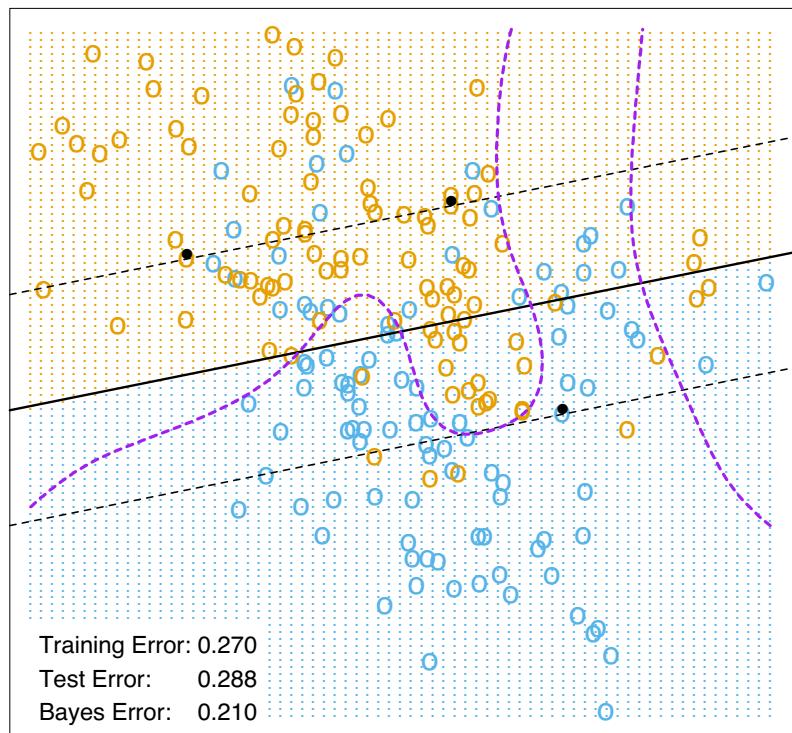
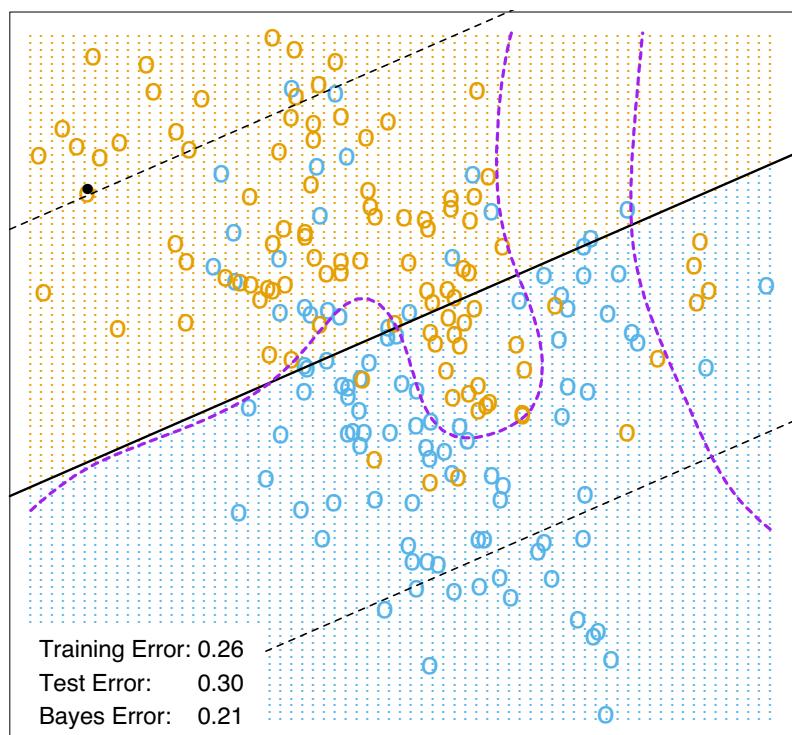
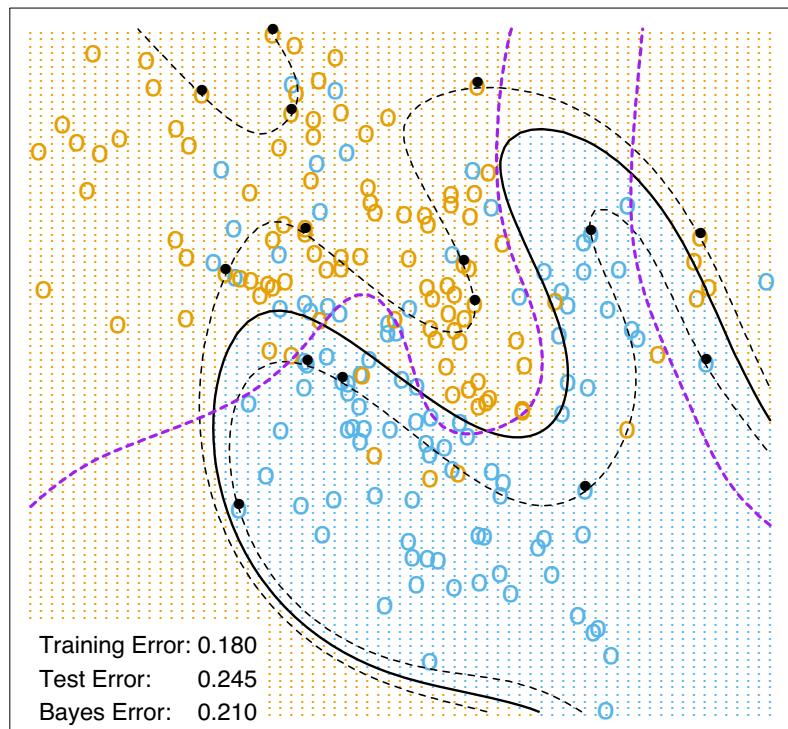

 $C = 10000$

 $C = 0.01$

FIGURE 12.2. The linear support vector boundary for the mixture data example with two overlapping classes, for two different values of C . The broken lines

SVM - Degree-4 Polynomial in Feature Space



SVM - Radial Kernel in Feature Space

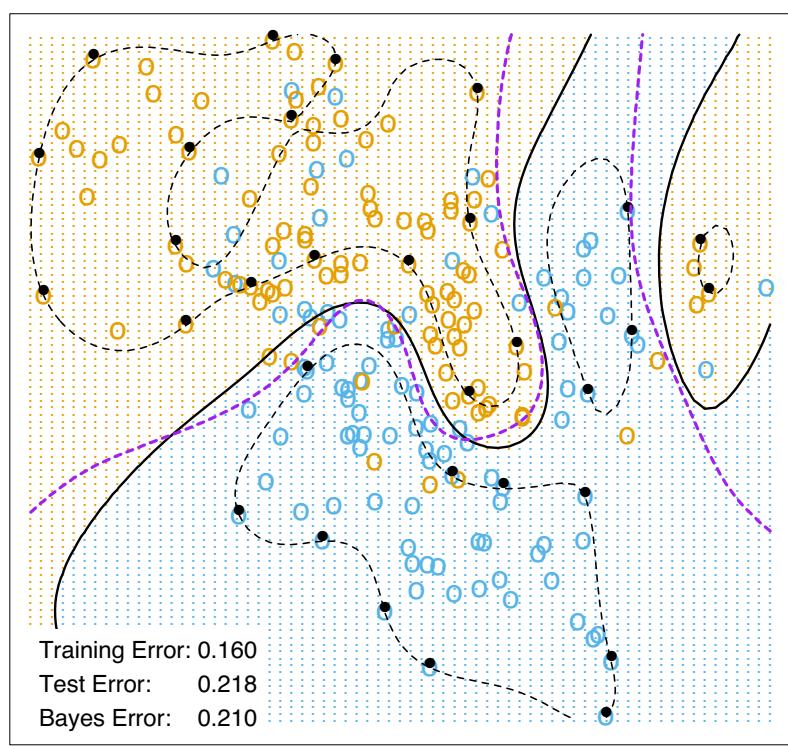


FIGURE 12.3. Two nonlinear SVMs for the mixture data. The upper plot uses a 4th degree polynomial

R K H S spaces of functions $\ni f$ classifier

infinite dim, with scalar product, complete

kernel

reproducing property (coming next)