

Lecture 12:

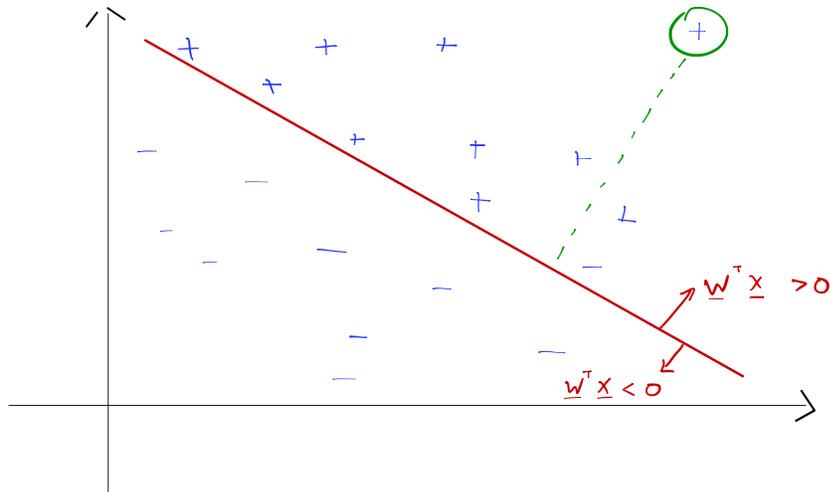
Support Vector Machines

# Support Vector Machines

Assume labels  $y_i \in \{+1, -1\}$

$$\text{Then } \|y - X\underline{w}\|_2^2 = \sum_{i=1}^n (y_i - \underbrace{x_i^T \underline{w}})^2 = \sum_{i=1}^n (1 - y_i \underbrace{x_i^T \underline{w}})^2$$

if this is  $\neq 1$ , then we incur loss even if point is correctly classified



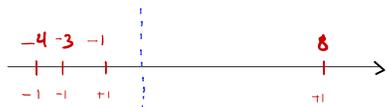
Ex. predict whether someone plays basketball based on height.

$y_i = +1 \rightarrow$  plays bball

$y_i = -1 \rightarrow$  doesn't play bball

4 training samples:  $x_i: 70', 71', 73', 82'$   $y_i: -1, -1, +1, +1$

$\xrightarrow{-4}$   $\xrightarrow{-3}$   $\xrightarrow{-1}$   $\xrightarrow{8}$   $\leftarrow$  difference from mean 74'

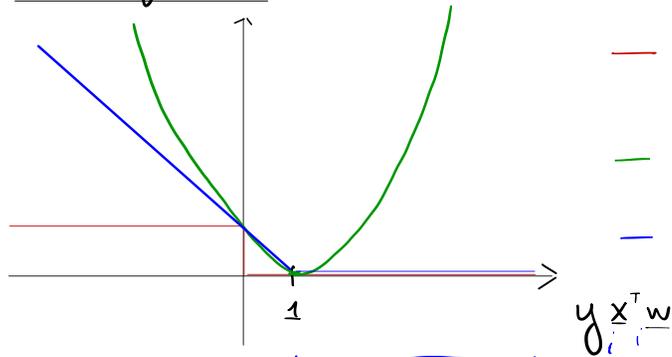


$$\hat{w}_{LS} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} = 0.15 \implies \hat{y} = +1 \text{ if } (\text{height} - \text{mean}) > 0.15$$

LS decision boundary

misclassifies a point even though perfect classifier exists

# Loss functions



some loss on some correctly classified points

no hinge loss, much quad loss

— 0/1 loss, =  $\mathbb{I}\{y \neq \text{sign}(x^T w)\} = \mathbb{I}\{y x^T w < 0\}$  = "ideal loss"

— quadratic loss =  $(1 - y x^T w)^2$

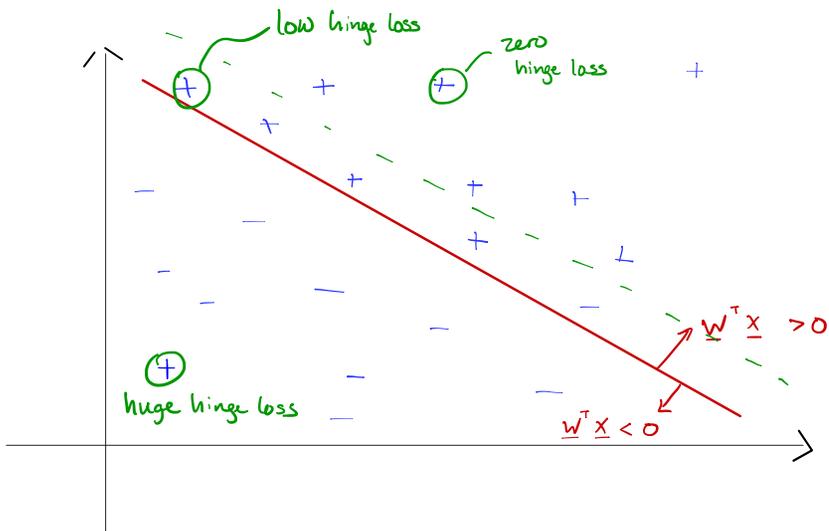
— hinge loss =  $(1 - y x^T w)_+$

Hinge loss mimics the ideal loss but is convex and easy to minimize

if  $y = +1, x^T w > 0$   
 then we accurately predicted label  
 if  $y = -1, x^T w < 0$ , acc label

$$(a)_+ = \begin{cases} a & \text{if } a > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{I}\{a\} = \begin{cases} 1 & \text{a true} \\ 0 & \text{a false} \end{cases}$$



$$l(\underline{w}) = \sum_{i=1}^n (1 - y_i \underline{x}_i^T \underline{w})_+$$

$$\nabla_{\underline{w}} l = \sum_{i=1}^n \mathbb{I}\{y_i \underline{x}_i^T \underline{w} < 1\} (-y_i \underline{x}_i)$$

↑ technically called "subgradient" because  $l(\underline{w})$  not differentiable

if  $y x^T w > 0$ , then correct label prediction  
 $y x^T w < 0$ , then wrong label prediction

If we minimize  $\sum_{i=1}^n (1 - y_i \underline{x}_i^T \underline{w})_+ + \lambda \|\underline{w}\|_2^2$   
 or the kernel version  $\sum_{i=1}^n (1 - y_i \phi(\underline{x}_i)^T \underline{w})_+ + \lambda \|\underline{w}\|_2^2$   
 this is called a support vector machine

$$\text{Let } \hat{w}_{svm} = \underset{w}{\text{argmin}} \sum_{i=1}^n (1 - y_i \underline{x}_i^T w)_+ + \lambda \|w\|_2^2$$

$$\text{or } \hat{w}_{svm} = \underset{w}{\text{argmin}} \sum_{i=1}^n (1 - y_i \phi(\underline{x}_i)^T w)_+ + \lambda \|w\|_2^2$$

Just like with kernel ridge regression, it is possible to show  $\hat{w}_{svm} = X^T \underline{\alpha}$  for some  $\underline{\alpha}$  (different from least squares  $\underline{\alpha}$ )

$$\text{or } \hat{w}_{svm} = \Phi^T \underline{\alpha}$$

To see this:

$$\text{Imagine } \hat{w} = X^T \underline{\alpha} + w^\perp = \sum_{j=1}^n \alpha_j \underline{x}_j + w^\perp \quad (w^\perp \text{ some vector orthogonal to the } \underline{x}_i\text{'s})$$

$$(a)_+ = \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{then } \min_{\underline{\alpha}, w^\perp} \sum_{i=1}^n (1 - y_i \underline{x}_i^T (\sum_{j=1}^n \alpha_j \underline{x}_j + w^\perp))_+ + \lambda \left\| \sum_{j=1}^n \alpha_j \underline{x}_j + w^\perp \right\|_2^2$$

$$= \min_{\underline{\alpha}, w^\perp} \sum_i (1 - y_i \left[ \sum_{j=1}^n \alpha_j \langle \underline{x}_i, \underline{x}_j \rangle + \underbrace{\underline{x}_i^T w^\perp}_{\text{always zero!}} \right])_+ + \lambda \left[ \left\| \sum_{j=1}^n \alpha_j \underline{x}_j \right\|_2^2 + \|w^\perp\|_2^2 \right]$$

$$= \min_{\underline{\alpha}, w^\perp} \sum_i (1 - y_i \left[ \sum_{j=1}^n \alpha_j \langle \underline{x}_i, \underline{x}_j \rangle \right])_+ + \lambda \left[ \left\| \sum_{j=1}^n \alpha_j \underline{x}_j \right\|_2^2 + \|w^\perp\|_2^2 \right]$$

$\rightarrow$  this must be 0 at optimum!

$$\hat{\underline{\alpha}} = \underset{\underline{\alpha}}{\text{argmin}} \sum_i (1 - y_i \left[ \sum_{j=1}^n \alpha_j \langle \underline{x}_i, \underline{x}_j \rangle \right])_+ + \lambda \left\| \sum_{j=1}^n \alpha_j \underline{x}_j \right\|_2^2$$

$$= \sum_i \sum_j \alpha_i \alpha_j \langle \underline{x}_i, \underline{x}_j \rangle$$

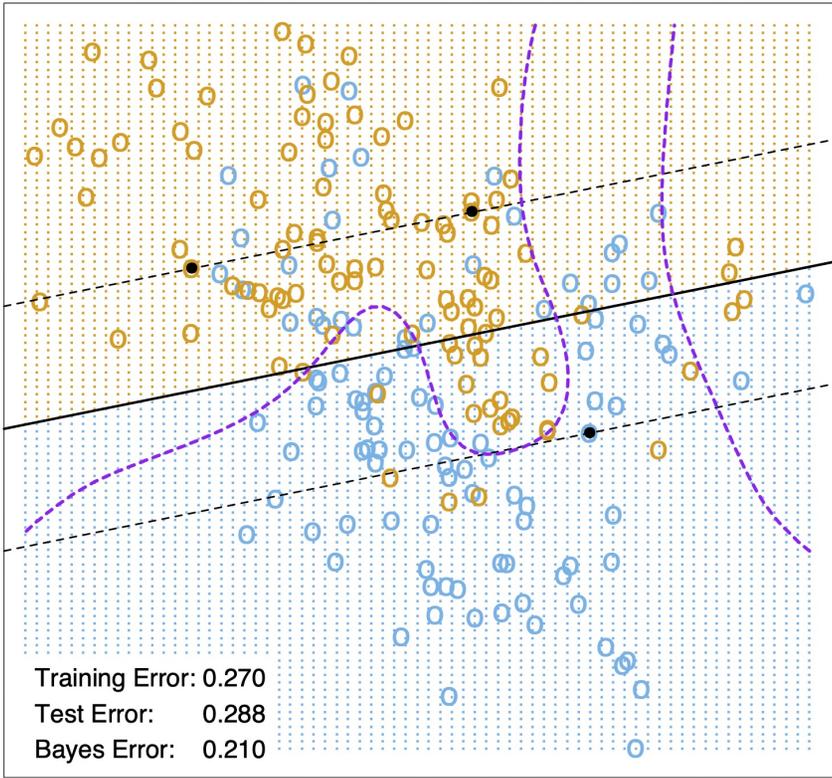
Alternatively, we can apply the "kernel trick" and replace inner products with kernels:

$$\hat{\underline{\alpha}} = \underset{\underline{\alpha}}{\text{argmin}} \sum_{i=1}^n (1 - y_i \sum_j \alpha_j K(\underline{x}_i, \underline{x}_j))_+ + \lambda \sum_i \sum_j \alpha_i \alpha_j K(\underline{x}_i, \underline{x}_j)$$

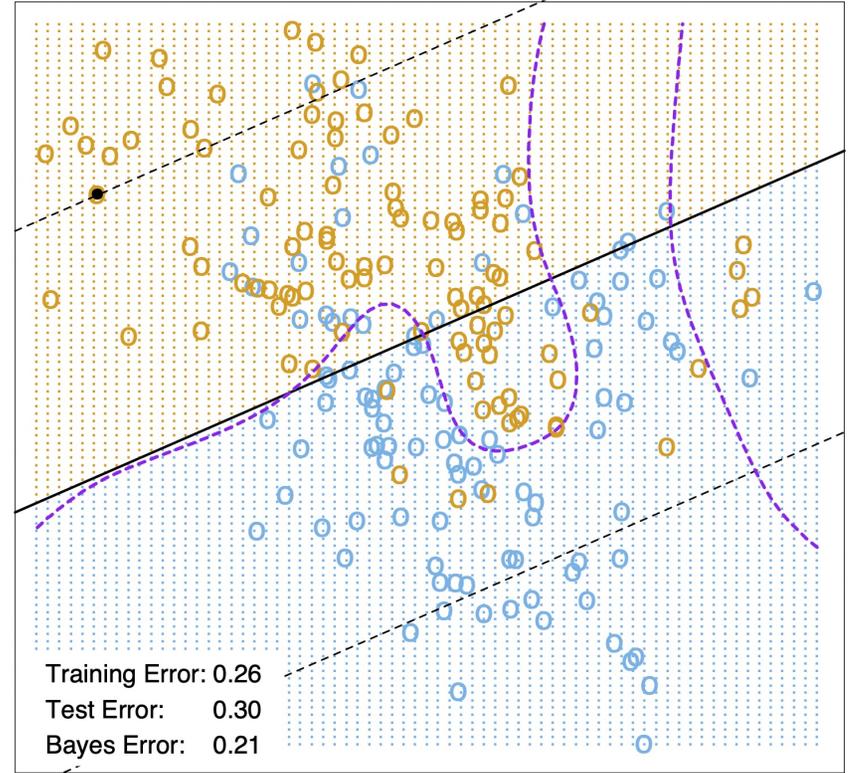
There is no closed form solution to this optimization problem.

$\Rightarrow$  need to use gradient descent or other numerical optimization methods.





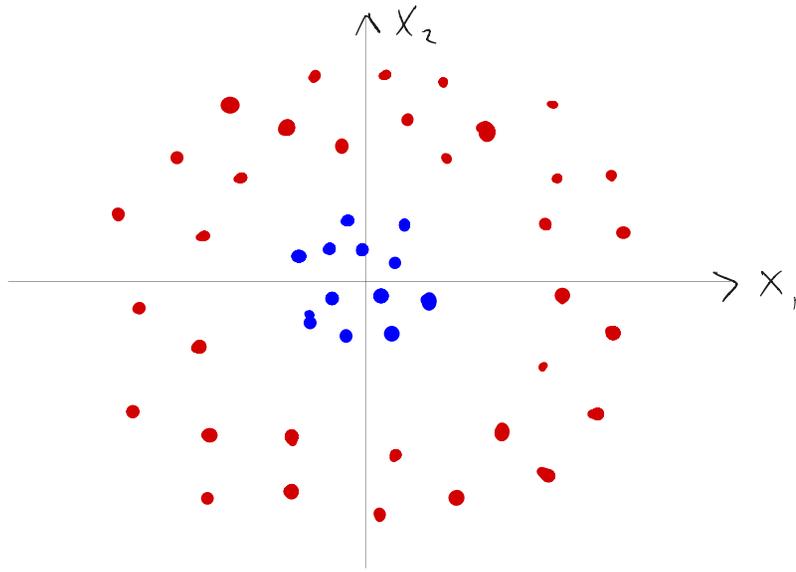
$$\lambda = 10^{-4}$$



$$\lambda = 100$$

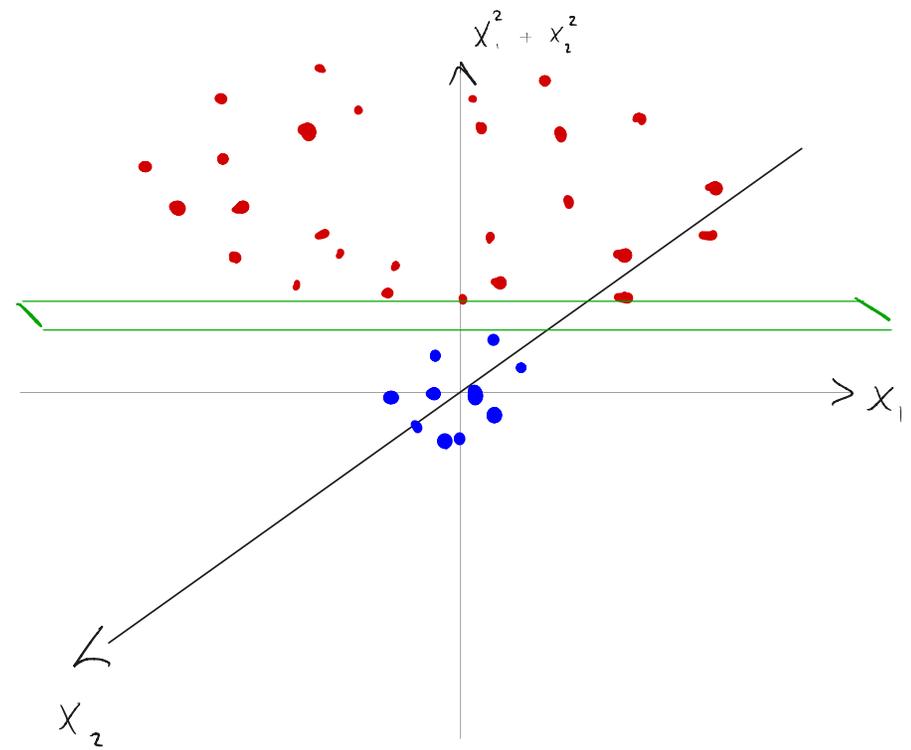
- ideal ('Bayes') decision boundary — gives smallest test errors on average
- points  $\underline{x}$  where  $\hat{y} = \underline{x}^T \hat{\underline{w}}_{svm} = \pm 1$  ←
- SVM decision boundary

# How do kernels help?



no good linear classifier exists

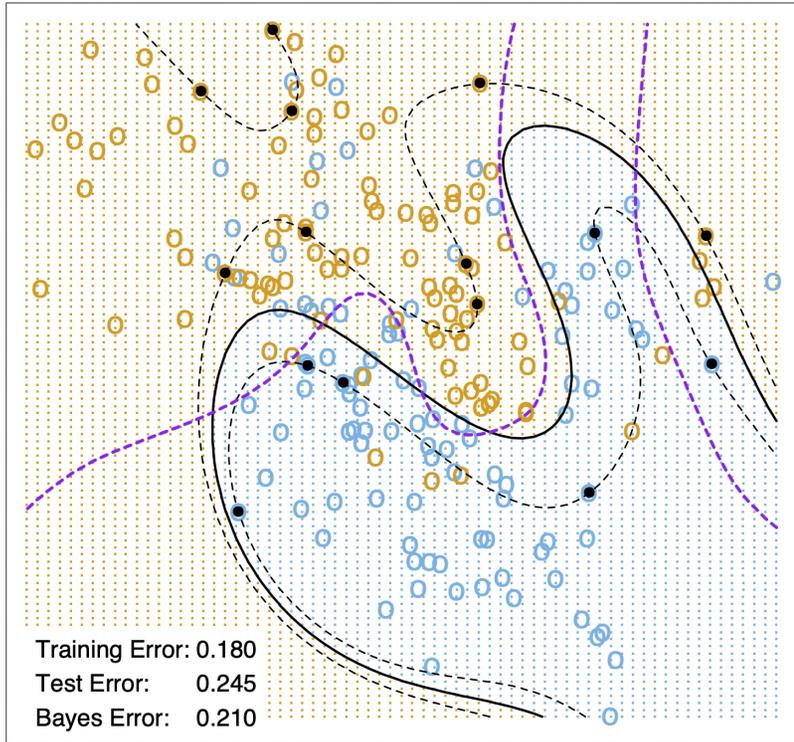
$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix}$$



in high-dimensional feature space, a good linear separating hyperplane exists.

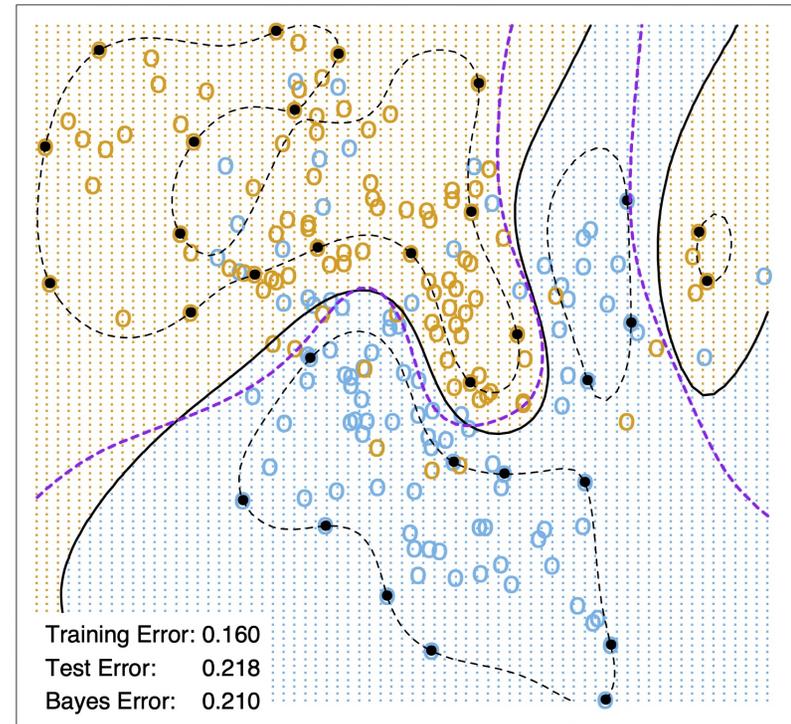
$$\phi(x_i) = \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i1}^2 + x_{i2}^2 \end{bmatrix}$$

SVM - Degree-4 Polynomial in Feature Space



(Gaussian)

SVM - Radial Kernel in Feature Space



- - - - - ideal ("Bayes") decision boundary — gives smallest test errors on average
- - - - - points  $x$  where  $\hat{y} = x^T \hat{w}_{svm} = \pm 1$
- SVM decision boundary

$$k(x_i, x_j) = \exp \left\{ - \frac{\|x_i - x_j\|_2^2}{2\sigma^2} \right\}$$