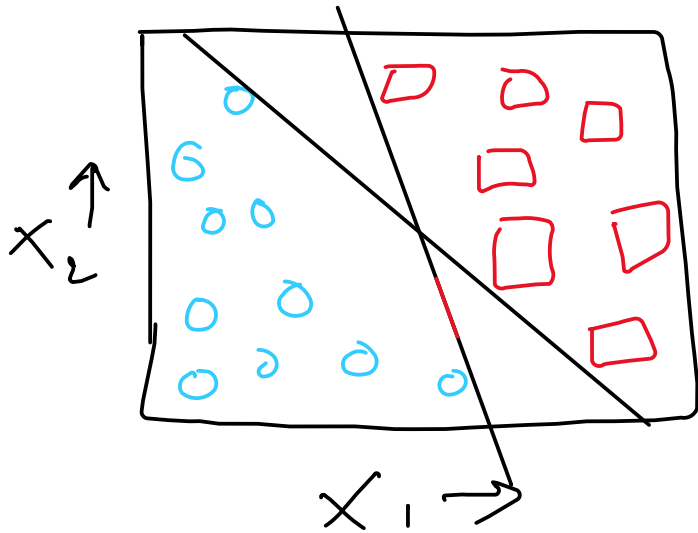


# Support vector machine

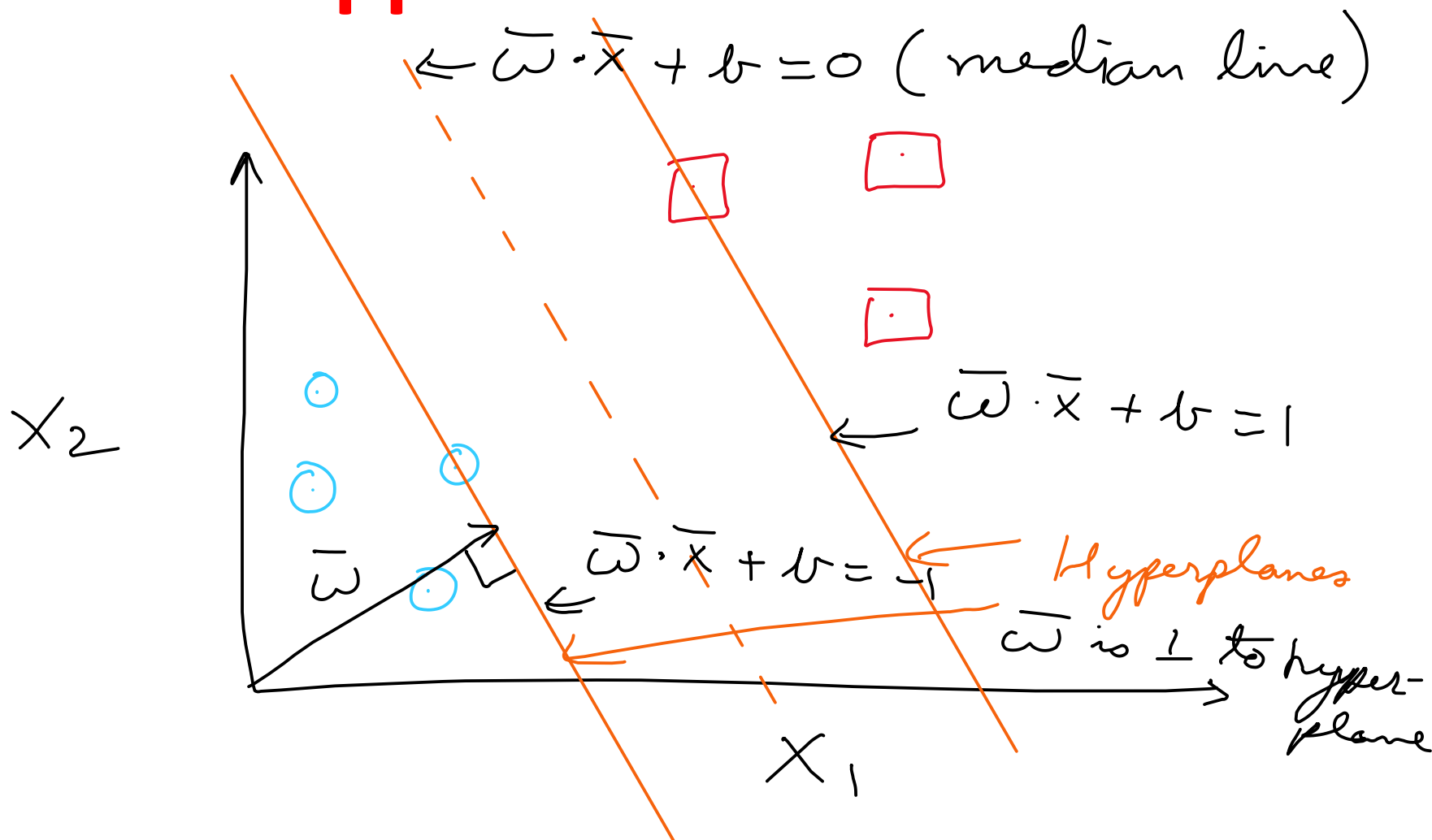
Prof. Asim Tewari  
IIT Bombay

# Support Vector Classifier

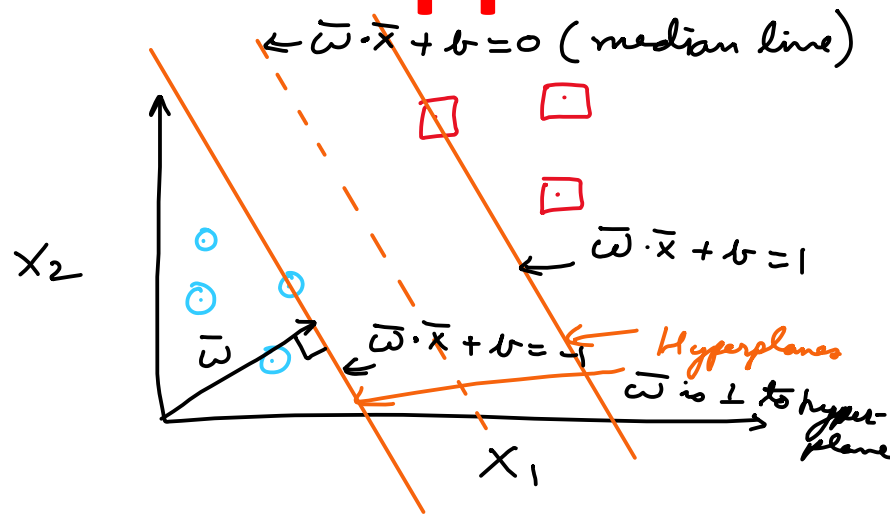


SVC : Maximum-margin  
hyperplane

# Support Vector Classifier



# Support Vector Classifier



① Any point above or on  $\bar{w} \cdot \bar{x} + b = 1$  belongs to  $\square$

② Any point on or below  $\bar{w} \cdot \bar{x} + b = -1$  belongs to  $\circ$

For " $\square$ ":  $\bar{w} \cdot \bar{x} + b \geq 1$

For " $\circ$ ":  $\bar{w} \cdot \bar{x} + b \leq -1$

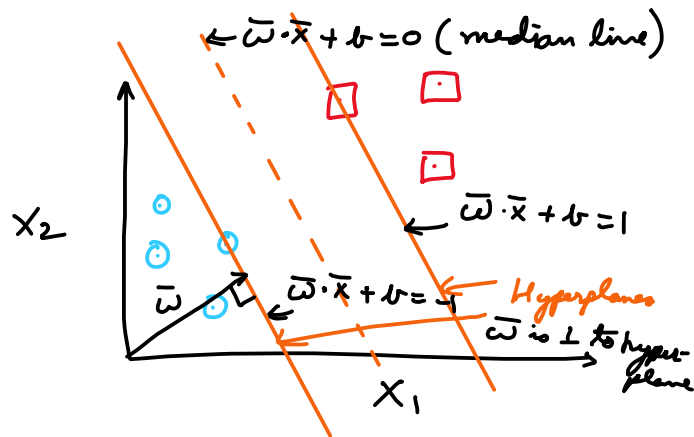
$\bar{w} \cdot \bar{x} + b \geq 0 \Rightarrow \square$

$\bar{w} \cdot \bar{x} + b \leq 0 \Rightarrow \circ$

For an unknown  $\bar{x}$ , the decision rule is based on maximum margin hyperplane.

# Support Vector Classifier

Standardize data, set  
For point  $i$



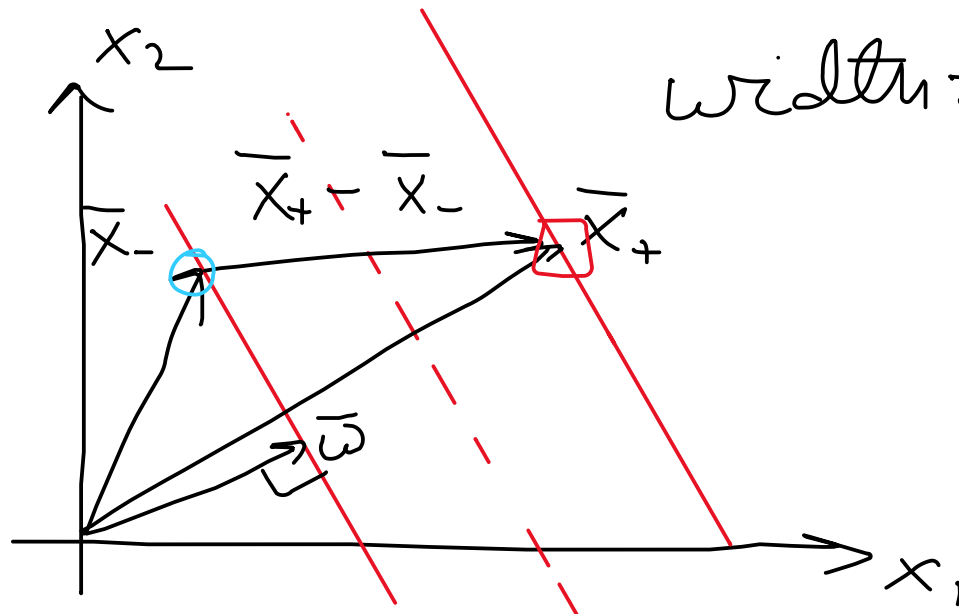
$$y_i = \begin{cases} 1 & \text{if it belongs to } \square \\ -1 & \text{if it belongs to } \circ \end{cases}$$

$$y_i (\bar{w} \cdot \bar{x}_i + b) \geq 1 \quad \left\{ \begin{array}{l} \square \text{ True for} \\ \circ \text{ both} \end{array} \right.$$

$$\Rightarrow y_i (\bar{w} \cdot \bar{x}_i + b) - 1 \geq 0$$

$$y_i (\bar{w} \cdot \bar{x}_i + b) - 1 = 0$$

# Support Vector Classifier



$$\text{width} = (\bar{x}_+ - \bar{x}_-) \cdot \frac{\bar{w}}{\|\bar{w}\|}$$

↑

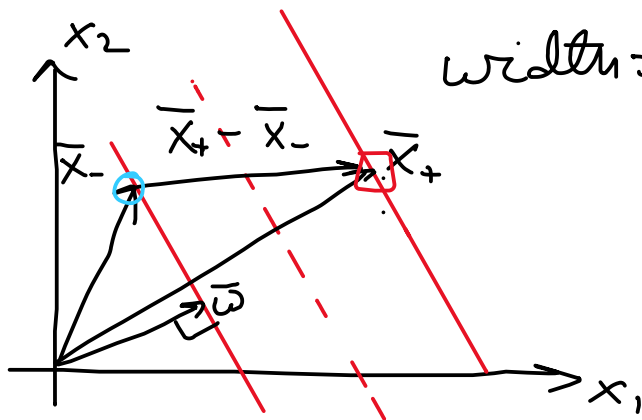
unit vector

$$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$$

$$\text{we get } \bar{w} \cdot \bar{x}_+ = 1 - b \text{ and } \bar{w} \cdot \bar{x}_- = -1 - b$$

$$\text{width} = (\bar{x}_+ - \bar{x}_-) \cdot \frac{\bar{w}}{\|\bar{w}\|} = \frac{\bar{x}_+ \cdot \bar{w} - \bar{x}_- \cdot \bar{w}}{\|\bar{w}\|}$$

# Support Vector Classifier



$$\text{width} = \frac{\bar{x}_+ \cdot \bar{w} - \bar{x}_- \cdot \bar{w}}{\|\bar{w}\|}$$

$$= \frac{(1 - \cancel{w}) - (-1 - \cancel{w})}{\|\bar{w}\|}$$

$$\Rightarrow \text{width} = \frac{2}{\|\bar{w}\|}$$

So we want two hyperplanes

$y_i(\bar{w} \cdot \bar{x}_i + b) - 1 = 0$  such that the width  $(\frac{2}{\|\bar{w}\|})$  is maximum.

Maximize  $\frac{2}{\|\bar{w}\|}$

$\Rightarrow$  Minimize  $\|\bar{w}\|/2$

# Support Vector Classifier

Minimize  $\|\omega\|_2$  with a constraint  $\equiv \min_{\omega} \|\omega\|^2$

$$y_i(\bar{\omega} \cdot \bar{x}_i + b) - 1 = 0$$

$$L = \frac{1}{2} \|\bar{\omega}\|^2 - \sum a_i [y_i(\bar{\omega} \cdot \bar{x}_i + b) - 1]$$

$$\min_{\omega, b}(L) \Rightarrow \frac{\partial L}{\partial \bar{\omega}} = 0, \quad \frac{\partial L}{\partial b} = 0$$

$$\Rightarrow \frac{\partial L}{\partial \bar{\omega}} = \frac{\partial}{\partial \bar{\omega}} \left[ \frac{1}{2} \|\bar{\omega}\|^2 \right] - \frac{\partial}{\partial \bar{\omega}} \left[ \sum a_i y_i (\bar{\omega} \cdot \bar{x}_i + b) - 1 \right]$$

$$\Rightarrow \bar{\omega} = \sum a_i y_i \bar{x}_i$$



# Support Vector Classifier

$$L = \frac{1}{2} \|\bar{\omega}\|^2 - \sum \alpha_i [y_i (\bar{\omega} \cdot x_i + b) - 1]$$

$$\frac{\partial L}{\partial \bar{\omega}} = 0 \Rightarrow \bar{\omega} = \sum (\alpha_i y_i \bar{x}_i)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow - \sum \alpha_i y_i = 0 \text{ or } \sum \alpha_i y_i = 0$$

$$\therefore L = \frac{1}{2} \underbrace{\left( \sum \alpha_i y_i \bar{x}_i \right)}_{\bar{\omega}} \cdot \underbrace{\left( \sum \alpha_j y_j \bar{x}_j \right)}_{\bar{\omega} = \|\bar{\omega}\|^2}$$

$$- \sum \left[ \alpha_i \left( y_i \left( \underbrace{\sum \alpha_j y_j x_j}_{\bar{\omega}} \right) \cdot \bar{x}_i + b \right) - 1 \right]$$

# Support Vector Classifier

$$L = \frac{1}{2} \left( \sum a_i y_i \bar{x}_i \right) \cdot \left( \sum a_j y_j \bar{x}_j \right) - \left( \sum a_i y_i \bar{x}_i \right) \cdot \left( \sum a_j y_j \bar{x}_j \right) - \sum a_i y_i b + \sum a_i$$

$$\Rightarrow L = \sum a_i - \frac{1}{2} \sum_i \sum_j \left( a_i a_j y_i y_j \bar{x}_i \cdot \bar{x}_j \right)$$

↑  
dot product

$$L = f(\bar{x}_i \cdot \bar{x}_j)$$

# Support Vector Classifier

$$\Rightarrow L = \sum a_i - \frac{1}{2} \sum_i \sum_j (a_i a_j y_i y_j \bar{x}_i \cdot \bar{x}_j)$$

$L = f(\bar{x}_i \cdot \bar{x}_j)$

↑  
dot product

Now if  $\sum a_i y_i \underbrace{\bar{x}_i \cdot \bar{\mu}}_{\bar{x}_i \cdot \bar{\mu}} + b \geq 0$

Decision Rule

also depends on the dot product.

# Maximal Margin Classifier

- **Hyperplane:** In a  $p$ -dimensional space, a *hyperplane* is a flat affine subspace of hyperplane dimension  $p - 1$

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0$$

- If  $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p > 0$

then this tells us that  $X$  lies to one side of the hyperplane.  
On the other hand, if

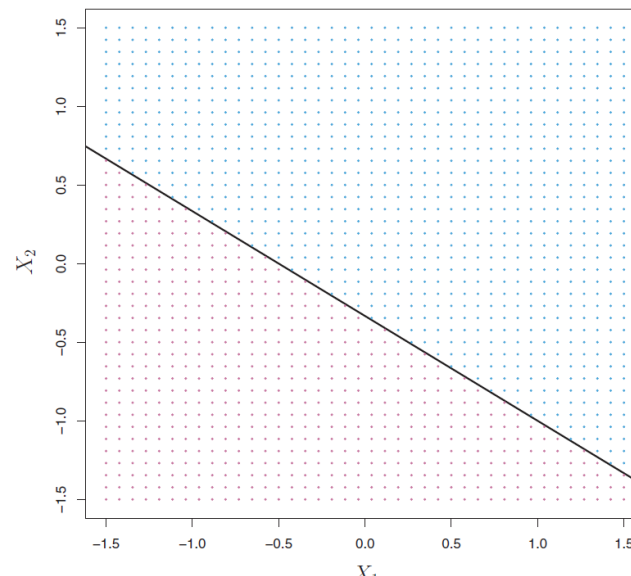
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p < 0$$

then  $X$  lies on the other side of the hyperplane. So we can think of the hyperplane as dividing  $p$ -dimensional space into two halves.

# Maximal Margin Classifier

- **Hyperplane:** In a  $p$ -dimensional space, a *hyperplane* is a flat affine subspace of hyperplane dimension  $p - 1$

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0$$



The hyperplane  $1 + 2X_1 + 3X_2 = 0$  is shown. The blue region is the set of points for which  $1 + 2X_1 + 3X_2 > 0$ , and the purple region is the set of points for which  $1 + 2X_1 + 3X_2 < 0$ .

# Maximal Margin Classifier

## *Classification Using a Separating Hyperplane*

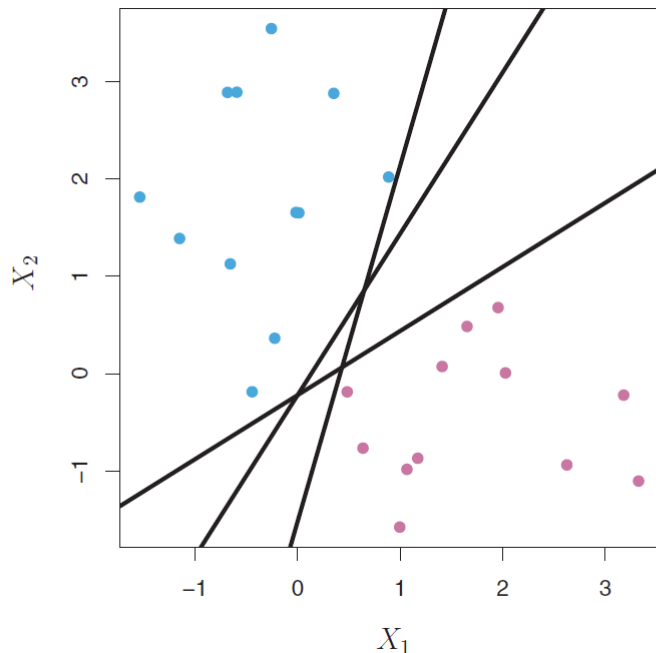
- Suppose that we have a  $n \times p$  data matrix  $X$ , that consists of  $n$  training observations in  $p$ -dimensional space,

$$x_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1p} \end{pmatrix}, \dots, x_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{np} \end{pmatrix}$$

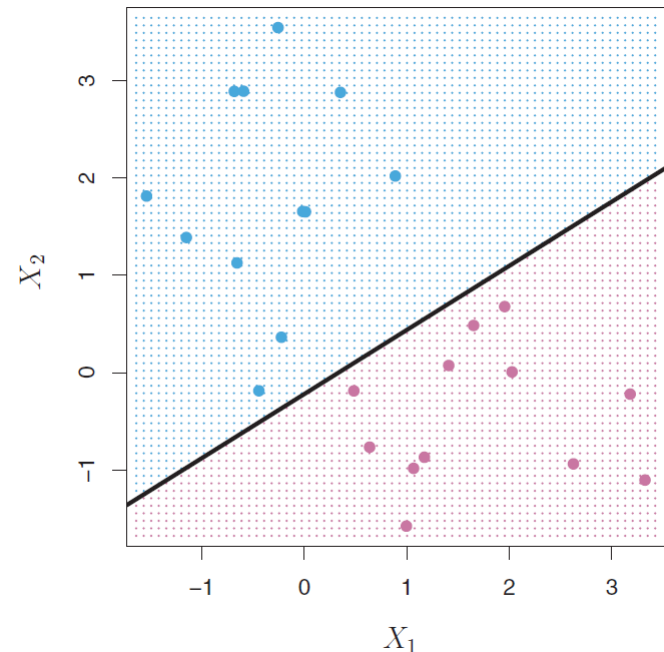
- and that these observations fall into two classes—that is,  $y_1, \dots, y_n \in \{-1, 1\}$  where  $-1$  represents one class and  $1$  the other class.

# Maximal Margin Classifier

## *Classification Using a Separating Hyperplane*



*There are two classes of observations, shown in blue and in purple, each of which has measurements on two variables. Three separating hyperplanes, out of many possible, are shown in black*



*A separating hyperplane is shown in black. The blue and purple grid indicates the decision rule made by Maximal Margin Classifier based on this separating hyperplane.*

# Maximal Margin Classifier

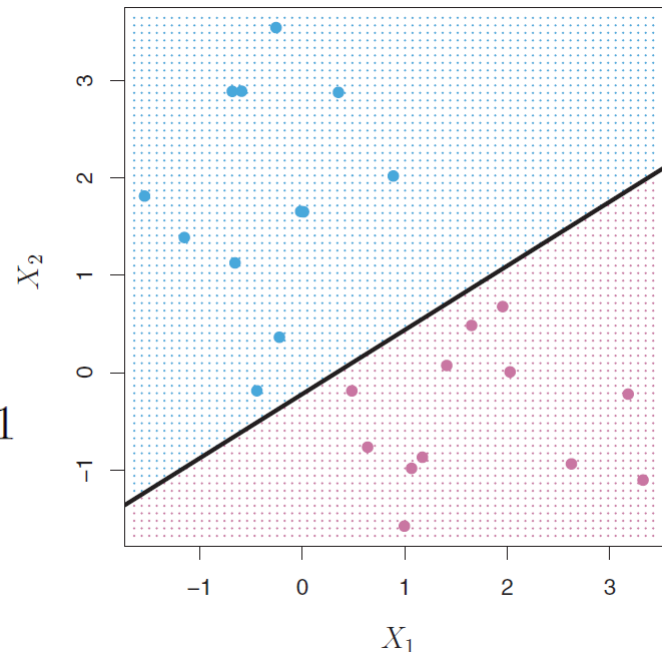
## *Classification Using a Separating Hyperplane*

We can label the observations from the blue class as  $y_i = 1$  and those from the purple class as  $y_i = -1$ . Then a separating hyperplane has the property that

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} > 0 \text{ if } y_i = 1$$

And

$$\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} < 0 \text{ if } y_i = -1$$



Equivalently, a separating hyperplane has the property that

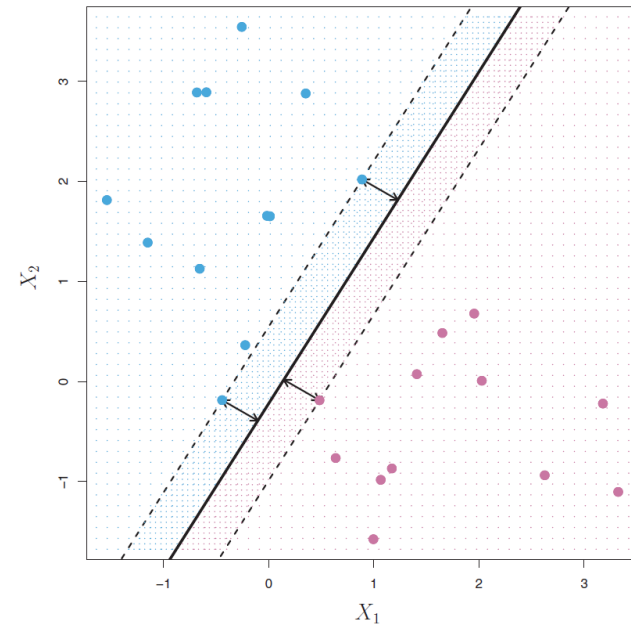
$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) > 0$$



# Maximal Margin Classifier

## *The Maximal Margin Classifier (separable case)*

*There are two classes of observations, shown in blue and in purple. The maximal margin hyperplane is shown as a solid line. The margin is the distance from the solid line to either of the dashed lines. The two blue points and the purple point that lie on the dashed lines are the support vectors, and the distance from those points to the hyperplane is indicated by arrows. The purple and blue grid indicates the decision rule made by a classifier based on this separating hyperplane.*



# Maximal Margin Classifier

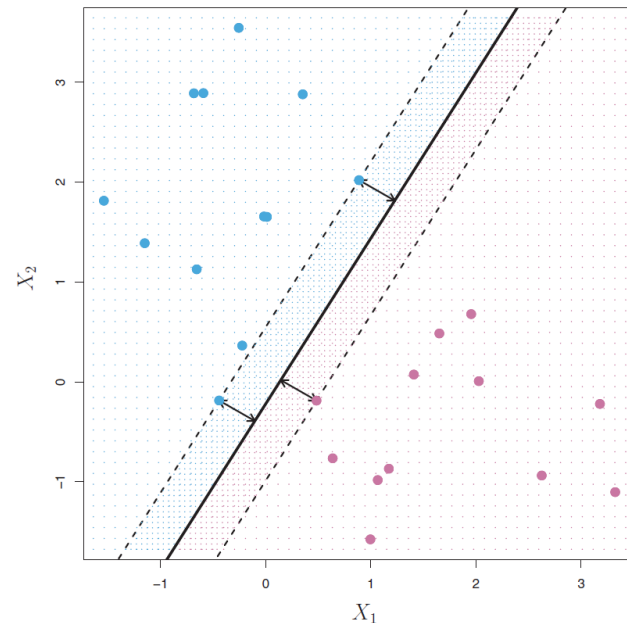
## *Construction of the Maximal Margin Classifier (separable Case)*

We now consider the task of constructing the maximal margin hyperplane based on a set of  $n$  training observations  $x_1, \dots, x_n \in \mathbb{R}_p$  and associated class labels  $y_1, \dots, y_n \in \{-1, 1\}$ . Briefly, the maximal margin hyperplane is the solution to the optimization problem

$$\underset{\beta_0, \beta_1, \dots, \beta_p}{\text{maximize}} \quad M$$

$$\text{subject to} \quad \sum_{j=1}^p \beta_j^2 = 1,$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n.$$



# Maximal Margin Classifier

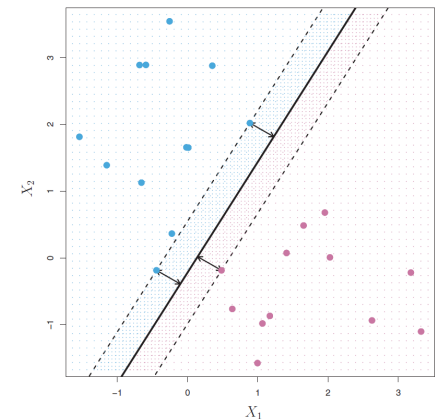
## *Construction of the Maximal Margin Classifier (separable Case)*

$$\text{maximize } M$$

$$\beta_0, \beta_1, \dots, \beta_p$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1,$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n.$$



- The above constraint guarantees that each observation will be on the correct side of the hyperplane, provided that  $M$  is positive.
- With this constraint the perpendicular distance from the  $i$ th observation to the hyperplane is given by  $y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip})$
- Therefore, the constraints ensure that each observation is on the correct side of the hyperplane and at least a distance  $M$  from the hyperplane.

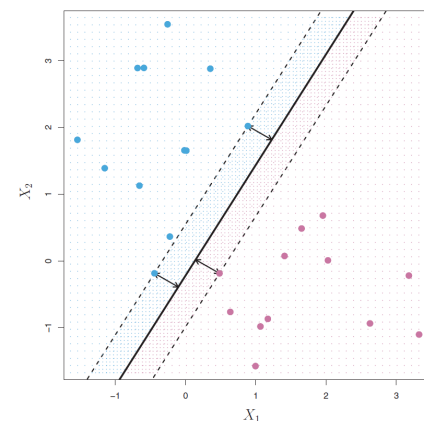
# Maximal Margin Classifier

## *Construction of the Maximal Margin Classifier (separable Case)*

maximize  $M$   
 $\beta_0, \beta_1, \dots, \beta_p$

subject to  $\sum_{j=1}^p \beta_j^2 = 1,$

$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n.$



- The above constrain guarantees that each observation will be on the correct side of the hyperplane, provided that  $M$  is positive.
- Actually, for each observation to be on the correct side of the hyperplane we would simply need  $y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) > 0$ , so the constraint in fact requires that each observation be on the correct side of the hyperplane, with some cushion, provided that  $M$  is positive.

# Maximal Margin Classifier

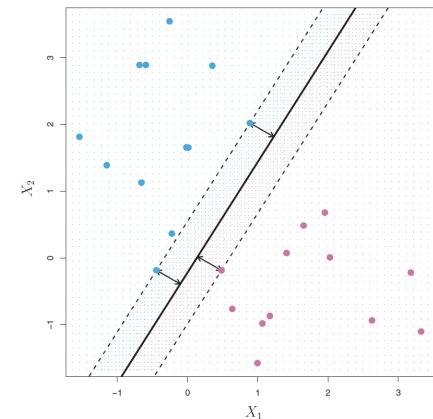
## *Construction of the Maximal Margin Classifier (separable Case)*

$$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1,$$

This is not really a constraint on the hyperplane, since if  $y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) = 0$  defines a hyperplane, then so does  $k y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) = 0$  for any  $k \neq 0$ . However, this constraint adds meaning to

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n.$$

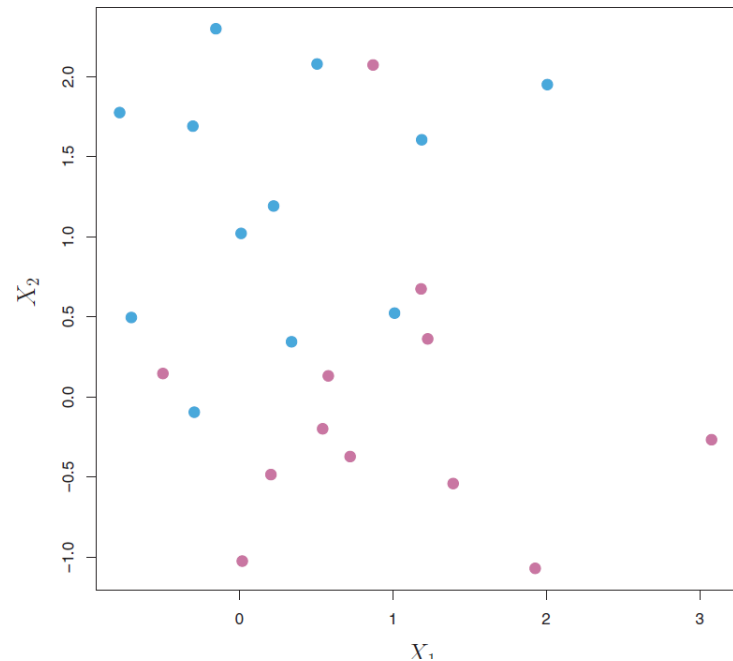
One can show that with this constraint the perpendicular distance from the  $i$ th observation to the hyperplane is given by  $y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip})$



# Support Vector Classifiers

## *The Maximal Margin Classifier (**Non-separable** Case)*

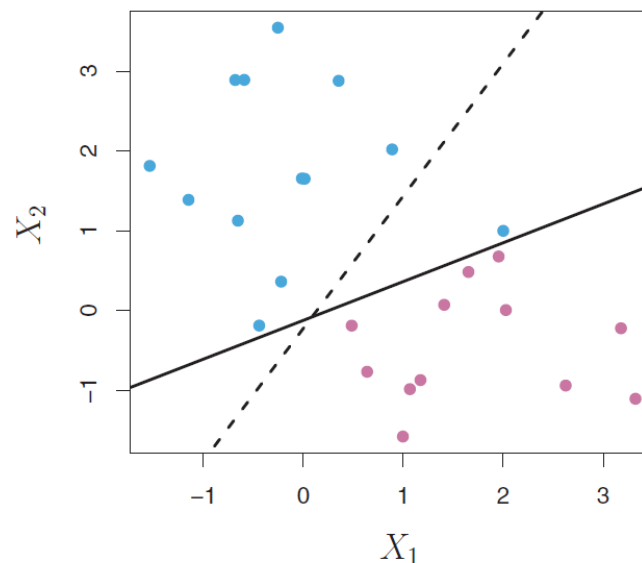
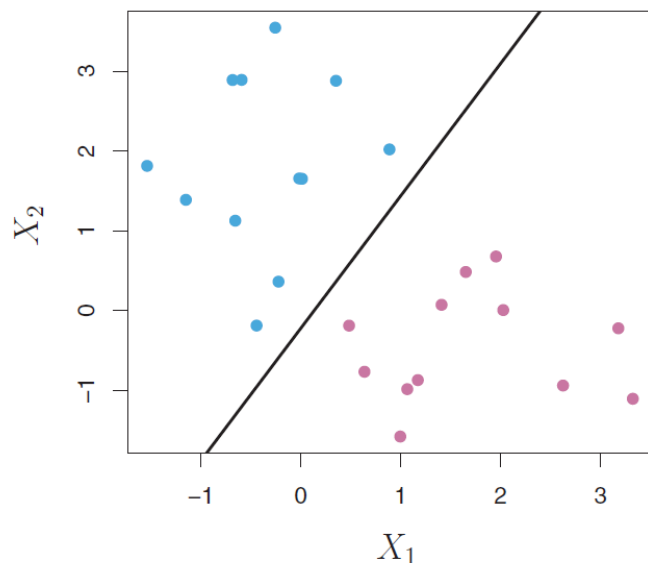
- In this case, the optimization problem has no solution with  $M > 0$ .
- Even if a separating hyperplane does exist, it would be very sensitive to individual observations.



*There are two classes of observations, shown in blue and in purple. In this case, the two classes are not separable by a hyperplane, and so the maximal margin classifier cannot be used.*

# Support Vector Classifiers

## *The Maximal Margin Classifier (**Non-separable** Case)*



The addition of a single observation could lead to a dramatic change in the maximal margin hyperplane.

# Support Vector Classifiers

## *The Maximal Margin Classifier (**Non-separable Case**)*

Thus, we want a classifier based on a hyperplane that does *not* perfectly separate the two classes, but has:

- Greater robustness to individual observations, and
- Better classification of *most* of the training observations.

The *support vector classifier*, sometimes called a *soft margin classifier*, does exactly this. Rather than seeking the largest possible margin so that every observation is not only on the correct side of the hyperplane but also on the correct side of the margin, we instead allow some observations to be on the incorrect side of the margin, or even the incorrect side of the hyperplane, to achieve **Greater robustness on *most* of the training observations.**



# Support Vector Classifiers

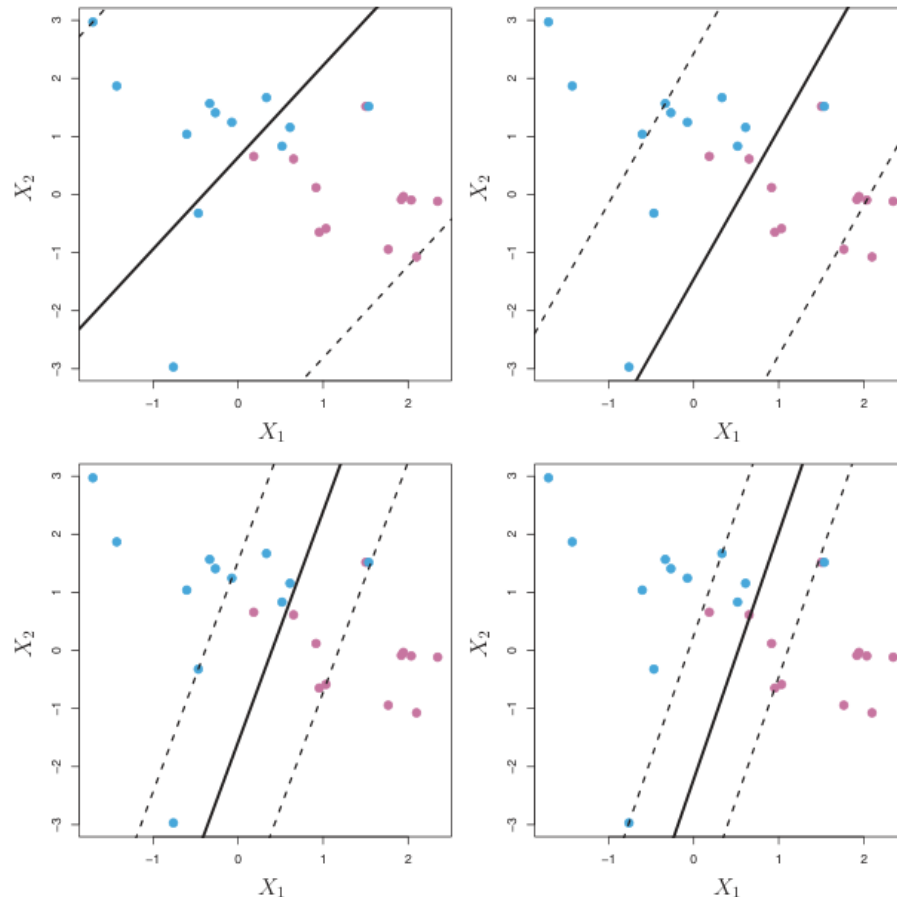
## *Construction of the Support Vector Classifier*

In this the hyperplane is chosen to correctly separate most of the training observations into the two classes, but may misclassify a few observations. It is the solution to the optimization problem

$$\begin{aligned} & \underset{\beta_0, \beta_1, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n}{\text{maximize}} && M \\ & \text{subject to} && \sum_{j=1}^p \beta_j^2 = 1, \\ & && y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i) \\ & && \epsilon_i \geq 0, \quad \sum_{i=1}^n \epsilon_i \leq C, \end{aligned}$$

where  $C$  is a nonnegative tuning parameter. We seek to make  $M$  as large as possible.  $\epsilon_1, \dots, \epsilon_n$  are *slack variables* that allow individual observations to be on the wrong side of the margin or the hyperplane.

# Support Vector Classifiers



A support vector classifier was fit using four different values of the tuning parameter  $C$ . The largest value of  $C$  was used in the top left panel, and smaller values were used in the top right, bottom left, and bottom right panels. When  $C$  is large, then there is a high tolerance for observations being on the wrong side of the margin, and so the margin will be large. As  $C$  decreases, the tolerance for observations being on the wrong side of the margin decreases, and the margin narrows.

# Mercer's theorem

Mercer's theorem states that for any data set  $X$  and any kernel function

$k : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$  there exists a mapping  $\varphi : \mathbb{R}^P \rightarrow \mathbb{R}^q$  so that

$$k(x_j, x_k) = \varphi(x_j) \cdot \varphi(x_k)^T$$

This means that a mapping from  $X$  to  $X'$  can be implicitly done by replacing scalar products in  $X'$  by kernel function values in  $X$ , without explicitly specifying  $\varphi$ , and without explicitly computing  $X'$ . Replacing scalar products in  $X'$  by kernel functions in  $X$  is called the *kernel trick*. Some frequently used kernel functions are

# Support Vector Classifier

$$\Rightarrow L = \sum a_i - \frac{1}{2} \sum_i \sum_j (a_i a_j y_i y_j \bar{x}_i \cdot \bar{x}_j)$$

$L = f(\bar{x}_i \cdot \bar{x}_j)$

↑  
dot product

Now if  $\sum a_i y_i \underbrace{\bar{x}_i \cdot \bar{\mu}}_{\bar{x}_i \cdot \bar{\mu}} + b \geq 0$

Decision Rule

also depends on the dot product.

# Support Vector Classifier

Decision rule as well as  $L$  depends only on the dot product.

$\therefore$  For a non-linear Hyper plane we can transform the original input space to another space by a

$$\bar{X} \xrightarrow{\Phi(\bar{X})} \bar{Z}$$

# Support Vector Classifier

$$\begin{array}{ccc} \bar{x} & \xrightarrow{\phi(\bar{x})} & \bar{z} \\ \text{You only need} & & \\ \bar{x}_i, \bar{x}_j & & \bar{z}_i, \bar{z}_j \end{array}$$

$$\bar{z}_i, \bar{z}_j$$

$$\Rightarrow \text{You only need to know } \phi(\bar{x}_i), \phi(\bar{x}_j)$$

From Mercer's theorem

$$\equiv \bar{z}_i, \bar{z}_j$$

This is achieved by a Kernel Function

$$K(\bar{x}_i, \bar{x}_j)$$

# Support Vector Classifier

Common Kernel Function

Polynomial  $k(\bar{x}_i, \bar{x}_j) = (\bar{x}_i \cdot \bar{x}_j)^a$   
(homogeneous)

(Non-Homogeneous)  $(\bar{x}_i \cdot \bar{x}_j + c)^a$

Gaussian radial basis function:

$$k(\bar{x}_i, \bar{x}_j) = \exp\left(\frac{-\|\bar{x}_i - \bar{x}_j\|^2}{2\sigma^2}\right)$$

# Support Vector Classifier

Hyperbolic tangent:

$$K(\bar{x}_i, \bar{x}_j) = \tanh(K \bar{x}_i \cdot \bar{x}_j + c)$$

for some  $K > 0, c < 0$

The decision boundary can be then given by

$$\bar{z}_i = \phi(\bar{x}_i), \quad \bar{z}_i \cdot \bar{z}_j = K(\bar{x}_i, \bar{x}_j)$$

$$\text{and } \bar{w} \cdot \bar{z}(i) = \bar{w} \cdot \phi(\bar{x}_i)$$

$$\Rightarrow \bar{w} \cdot \bar{z}_i = \sum a_j y_j \cdot K(\bar{x}_i, \bar{x}_j)$$



# The Support Vector Machine Kernel Trick

linear kernel

$$k(x_j, x_k) = x_j \cdot x_k^T$$

polynomial kernel

$$k(x_j, x_k) = (x_j \cdot x_k^T)^d, \quad d \in \{2, 3, \dots\}$$

Gaussian kernel

$$k(x_j, x_k) = e^{-\frac{\|x_j - x_k\|^2}{\sigma^2}}, \quad \sigma > 0$$

hyperbolic tangent kernel

$$k(x_j, x_k) = 1 - \tanh \frac{\|x_j - x_k\|^2}{\sigma^2}, \quad \sigma > 0$$

radial basis function (RBF) kernel

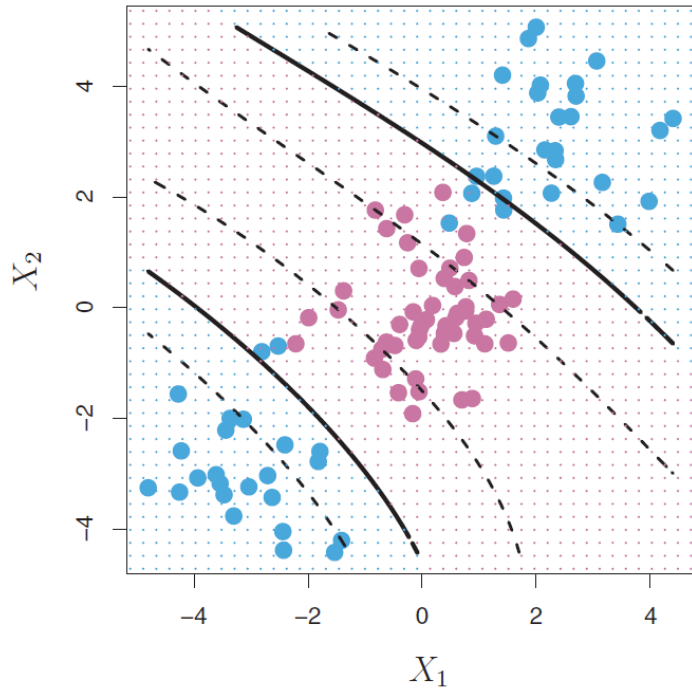
$$k(x_j, x_k) = f(\|x_j - x_k\|)$$

The SVM classification constraints are

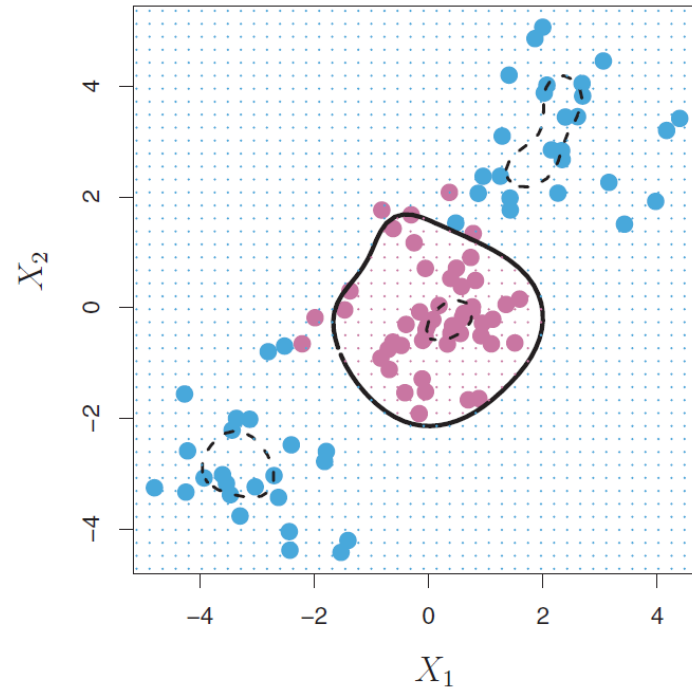
$$\sum_{y_j=1} \alpha_j k(x_j, x_k) - \sum_{y_j=2} \alpha_j k(x_j, x_k) + b \geq +1 - \xi_k \quad \text{if } y_k = 1$$

$$\sum_{y_j=1} \alpha_j k(x_j, x_k) - \sum_{y_j=2} \alpha_j k(x_j, x_k) + b \leq -1 + \xi_k \quad \text{if } y_k = 2$$

# The Support Vector Machine



*An SVM with a polynomial kernel of degree 3 is applied to the non-linear data, resulting in a far more appropriate decision rule*



*An SVM with a radial kernel is applied. In this example, either kernel is capable of capturing the decision boundary.*