

# Support Vector Machines

Roberta De Vito



**BROWN**  
Public Health

# Support vector machine

- ▶ Maximal Marginal Classifier
- ▶ It cannot be applied to most data sets
- ▶ Separation by a linear boundary

# Hyperplane

- ▶ Mathematical Definition

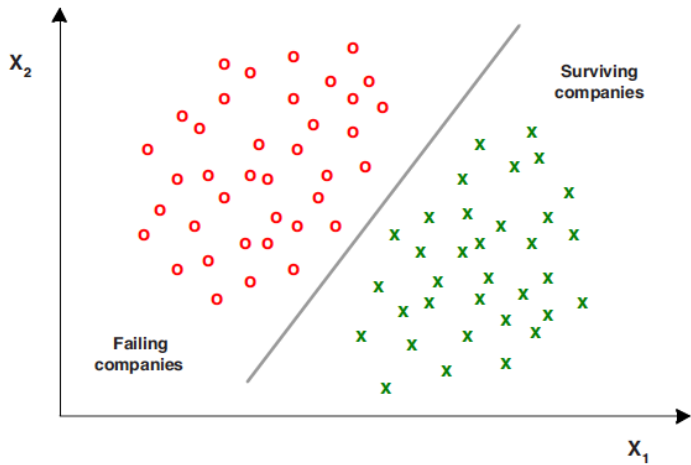
$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p = 0$$

- ▶ if  $X$  satisfies instead

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p > 0$$

What does it mean (Q1)?

# Hyperplane

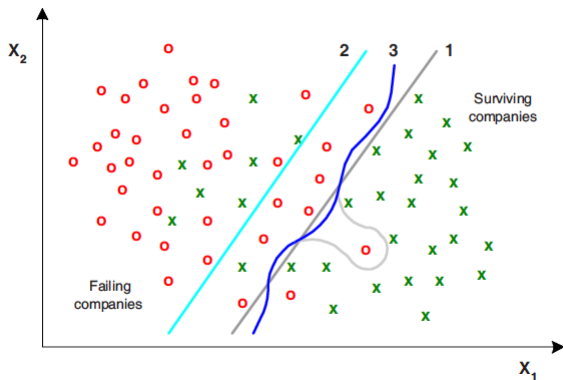


# Classification Using a Separating Hyperplane

- ▶ Our data matrix
- ▶ observations fall into two classes  $y_i \in \{-1, 1\}$

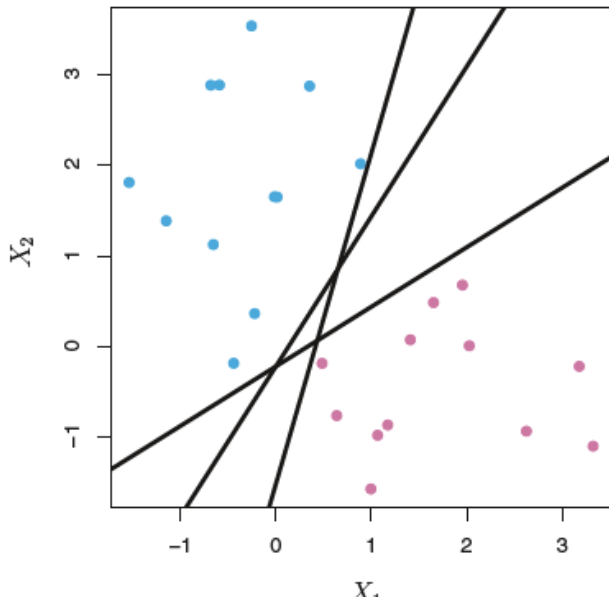
# How can we find the best separating hyperplanes

1 It should be linear



# How can we find the best separating hyperplanes

Also linear we can find many hyperplanes



# How can we find the best separating hyperplanes

A separating hyperplane has the property that

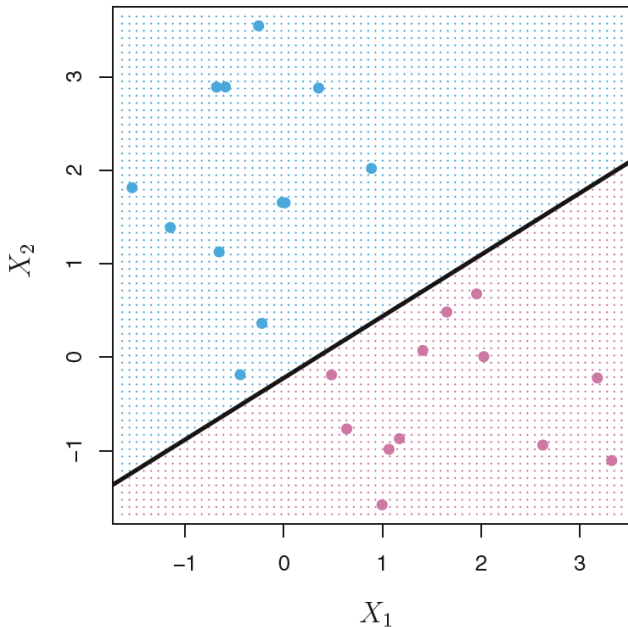
$$\begin{aligned}\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} &> 0 \text{ if } y_i = 1, \\ \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} &< 0 \text{ if } y_i = -1,\end{aligned}$$

Also we can write

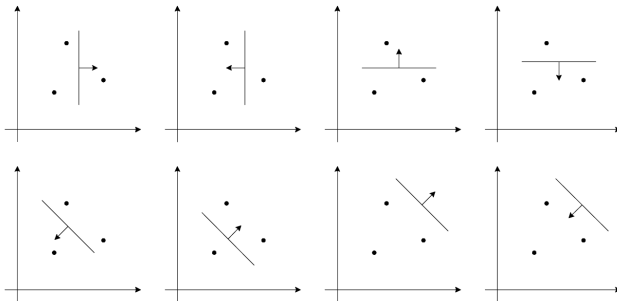
$$y_i(\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}) > 0$$



# How can we find the best separating hyperplanes



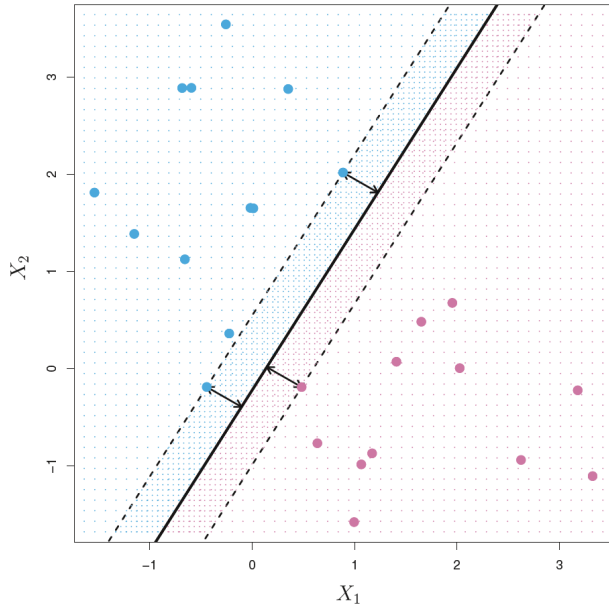
## 8 Different possible ways to divide three points



# The Maximal Margin Classifier

- ▶ look for the optimal separating hyperplane
- ▶ separating hyperplane that is farthest from the training observation
- ▶ perpendicular distance
- ▶ then, classify a test observation based the maximal margin hyperplane where it lies

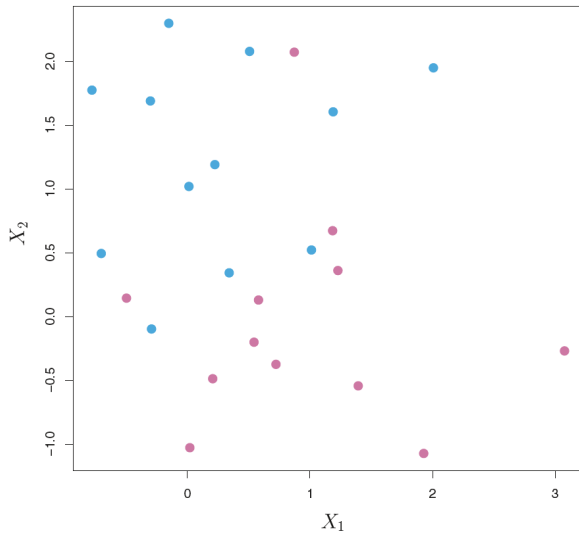
# The Maximal Margin Classifier



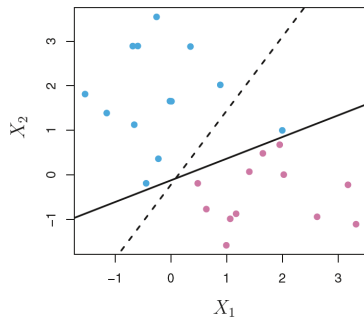
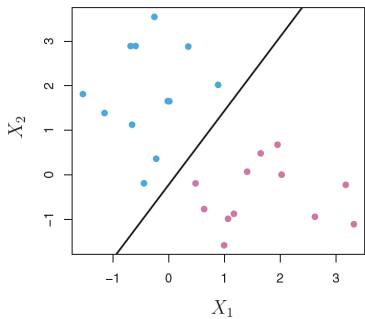
# The Maximal Margin Classifier

$$\begin{aligned} & \text{maximize } M \\ & \text{subject to } \sum_{i=1}^p \beta_j^2 = 1 \\ & y_i(\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}) \geq M \end{aligned}$$

# Observations not separable by a hyperplane



# Robustness (Q2)



# The soft margin classifier (Q3)

- ▶ margin is soft
- ▶ separate most of the training observations into the two classes

*maximize*  $M$

*subject to*  $\sum_{i=1}^p \beta_j^2 = 1$

$y_i(\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip}) \geq M(1 - \epsilon_i)$

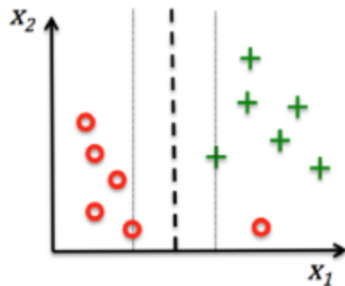
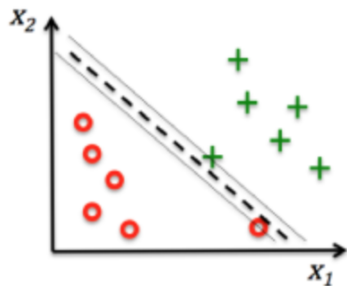
$\epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C$



# The severity of violations

- ▶  $C$  determines the severity of the violations to the margin
- ▶  $C = 0$ , no margin for violations  $\epsilon_1 = \dots = \epsilon_n = 0$
- ▶  $C > 0$  no more than  $C$  observations can be on the wrong side of the hyperplane
- ▶  $C$  is chosen via CV
- ▶ low-bias tradeoff (Q4)

## The severity of violations Q5



# Non-linear Decision Boundaries

- ▶  $2p$  features  $X_1, X_1^2, \dots, X_p, X_p^2$
- ▶ Now our equations will become

*maximize*  $M$

$$\sum_{j=1}^p \sum_{k=1}^2 \beta_{jk}^2 = 1$$

$$\text{subject to } y_i(\beta_0 + \sum_{j=1}^p \beta_{j1}x_{ij} + \sum_{j=1}^p \beta_{j2}x_{ij}^2) \geq M(1 - \epsilon_i)$$

$$\sum_{i=1}^n \epsilon_i \leq C, \epsilon_i \geq 0$$

# SVM with kernels

- ▶ Solution of the equations involves inner products of the observations

$$\langle x_i, x_{i'} \rangle = \sum_{j=1}^p x_{ij} x_{i'j}$$

- ▶ the linear support vector classifier

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$$

- ▶ kernel

$$K(x_i, x_{i'}) = \sum_{j=1}^p x_{ij} x_{i'j}$$

- ▶ the non-linear function now is

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$$

## SVM with kernels (Q6)

