

# Lecture 13. Support Vector Machines

COMP90051 Statistical Machine Learning

Semester 1, 2021  
Trevor Cohn



# This lecture

- Support vector machines (SVMs) as maximum-margin classifiers
- The hard-margin SVM objective
- SVM objective as regularised loss function
- The soft-margin SVM

# Maximum-Margin Classifier: Motivation

A new twist to binary linear classification

# Beginning: linear SVMs

- In the first part, we will consider a basic setup of SVMs, something called linear *hard-margin SVM*.
- Keep in mind: SVMs are more powerful than they may initially appear
- For now, we model the data as linearly separable, i.e., there exists a hyperplane perfectly separating the classes

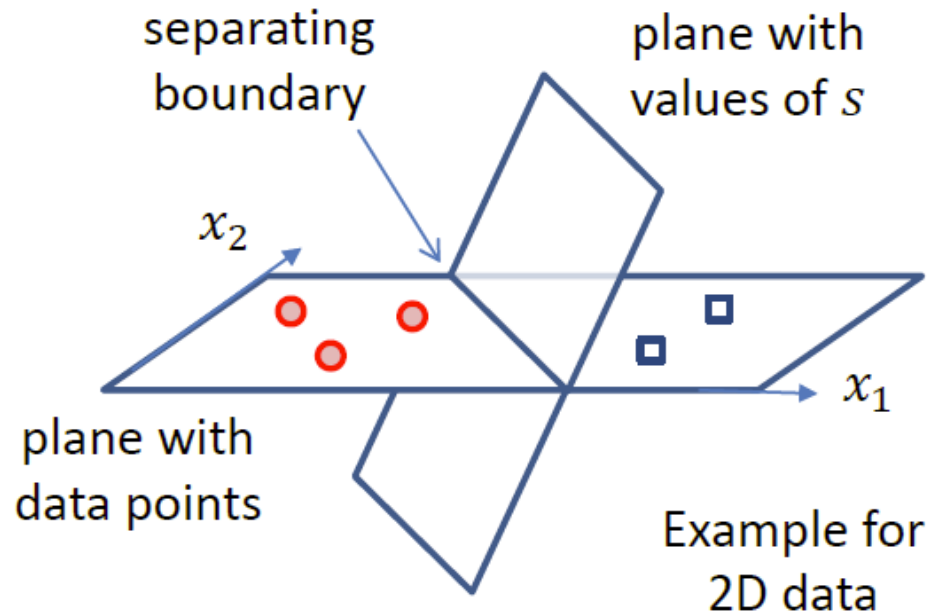
# SVM is a linear binary classifier

Predict class A if  $s \geq 0$

Predict class B if  $s < 0$

where  $s = b + \sum_{i=1}^m x_i w_i$

SVM is a linear classifier:  $s$  is a linear function of inputs, and the separating boundary is linear

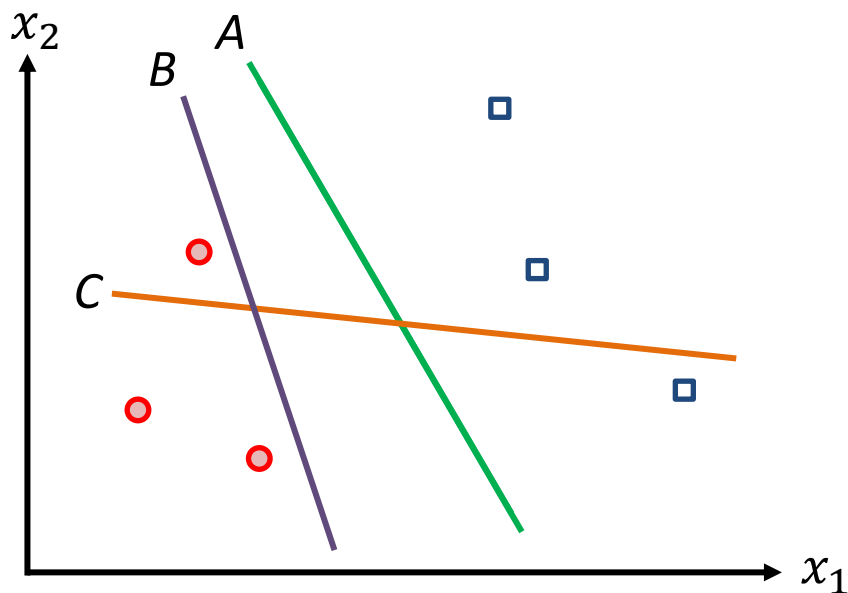


# SVM vs. perceptron, log. regression

- Exact same formulation as perceptron
- Same formulation for log-odds as log. regression
- How SVM is different: way parameters are learned.
  - \* Perceptron: min perceptron loss as studied earlier
  - \* Log. reg.: min binary cross-entropy loss
  - \* SVMs: different criterion for choosing parameters

# Choosing separation boundary

- An SVM is a linear binary classifier: choosing parameters means choosing a separating boundary (hyperplane)
- In 2D:

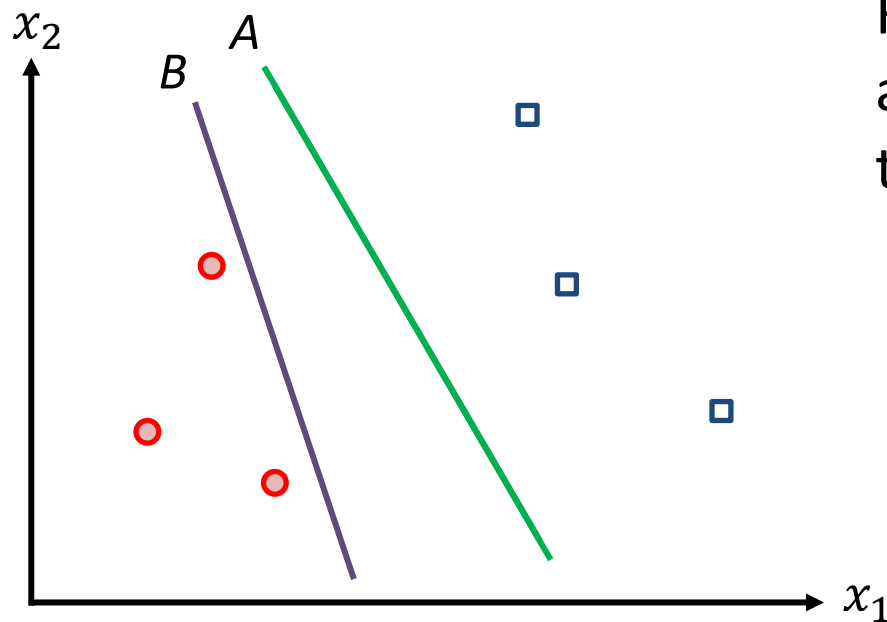


Which boundary would you choose?

A (Green)  
B (Purple)  
C (Orange)

# Which boundary should we use?

- Provided the dataset is linearly separable, the perceptron will find a boundary that separates classes perfectly. This can be any such boundary, e.g., A or B

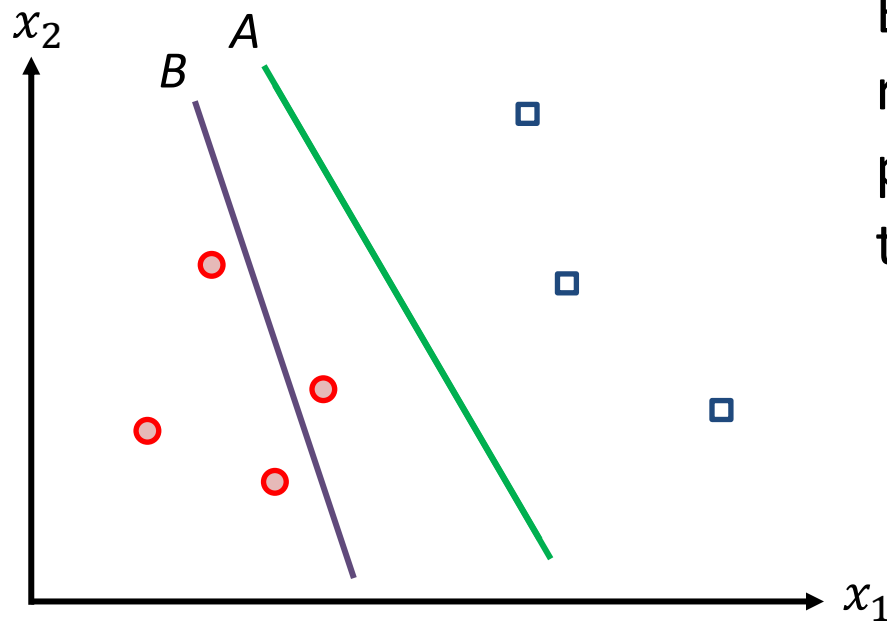


For the perceptron, they are equally good, because the perceptron loss is 0.



# Which boundary should we use?

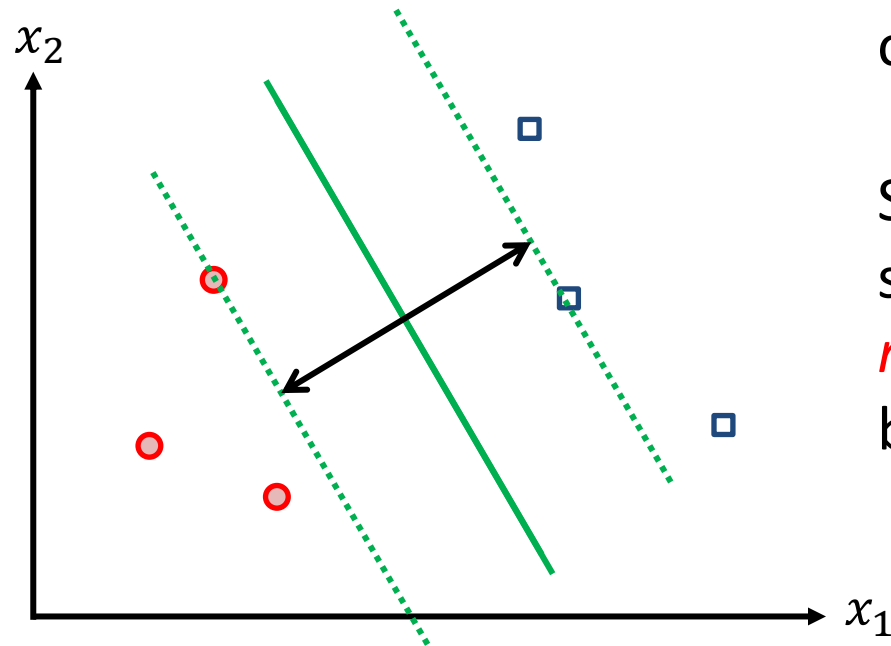
- Provided the dataset is linearly separable, the perceptron will find a boundary that separates classes perfectly. This can be any such boundary, e.g., A or B



But... line A seems more reliable. When new data point arrives, line B is likely to misclassify it

# Aiming for the safest boundary

- Intuitively, the most reliable boundary would be the one that is between the classes and as far away from both classes as possible



SVM objective captures this observation

SVMs aim to find the separation boundary that *maximises the margin* between the classes

# Maximum-margin classifier

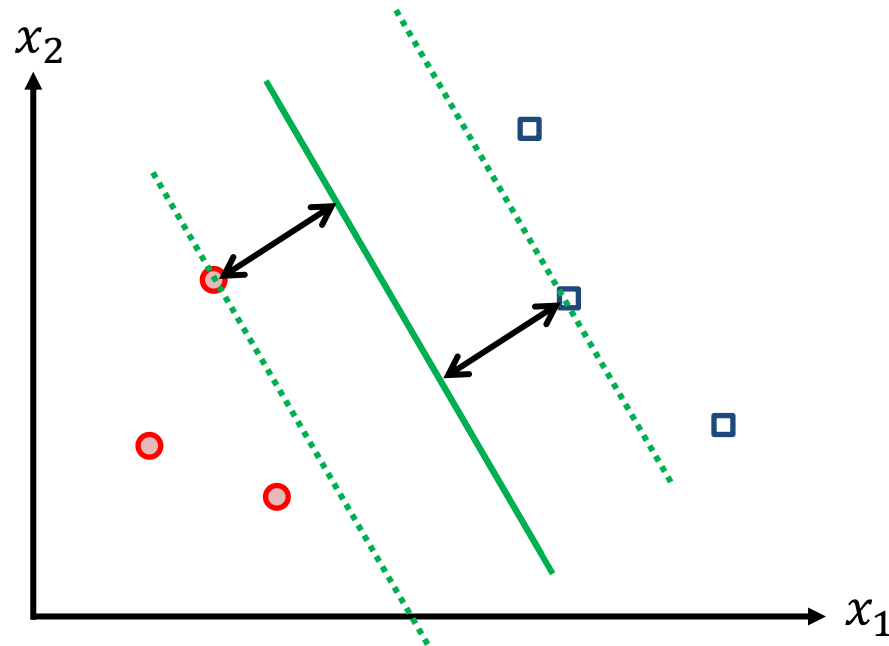
- An SVM is a linear binary classifier. SVM training aims to find the separating boundary that maximises margin
- For this reason, SVMs a.k.a *maximum-margin classifiers*
- The training data is fixed, so the margin is defined by the location and orientation of the separating boundary
- Our next step is to formalise our objective by expressing *margin width* as a function of parameters (and data)

# Maximum-Margin Classifier: Derivation

A geometric derivation of  
the SVM's objective

# Margin width

- While the margin can be thought as the space between two dashed lines, it is more convenient to define **margin width** as the distance between the separating boundary and the nearest data point(s)

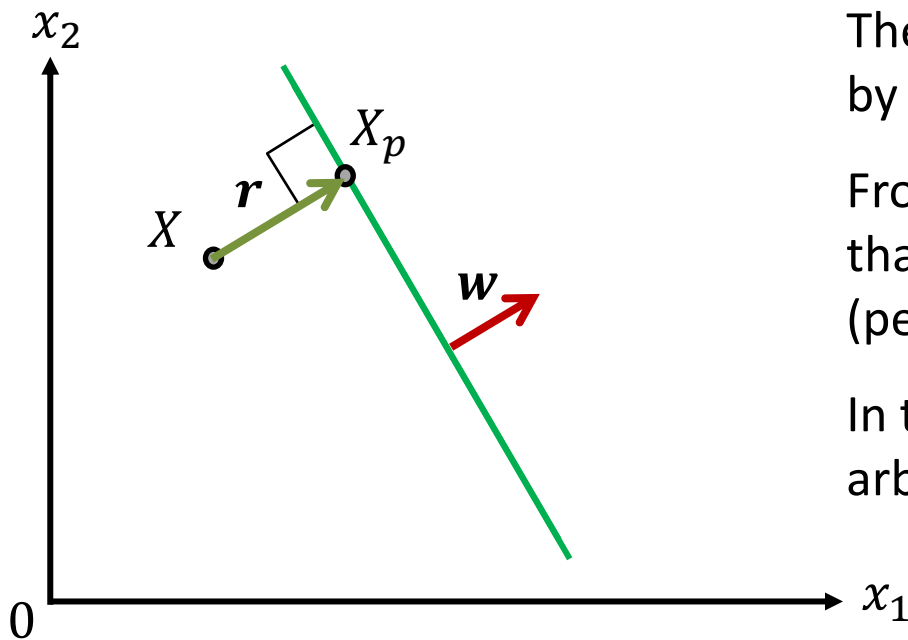


Point(s) on margin boundaries called *support vectors*

We want to maximise the distance to support vectors

# Distance from point to hyperplane

- Consider an arbitrary point  $X$  (from either of the classes, and not necessarily the closest one to the boundary), and let  $X_p$  denote the **projection** of  $X$  onto the separating boundary
- Now, let  $\mathbf{r}$  be a vector  $X_p - X$ . Note that  $\mathbf{r}$  is **perpendicular** to the boundary, and also that  $\|\mathbf{r}\|$  is the **required distance**



The separation boundary is defined by parameters  $\mathbf{w}$  and  $b$ .

From our linear algebra slides, recall that  $\mathbf{w}$  is a vector normal (perpendicular) to the boundary

In the figure,  $\mathbf{w}$  is drawn from an arbitrary starting point on boundary

# Distance from point to hyperplane

- Distance is  $\|\mathbf{r}\| = -\frac{\mathbf{w}'\mathbf{x}+b}{\|\mathbf{w}\|}$ , or more generally  $\|\mathbf{r}\| = \pm \frac{\mathbf{w}'\mathbf{x}+b}{\|\mathbf{w}\|}$

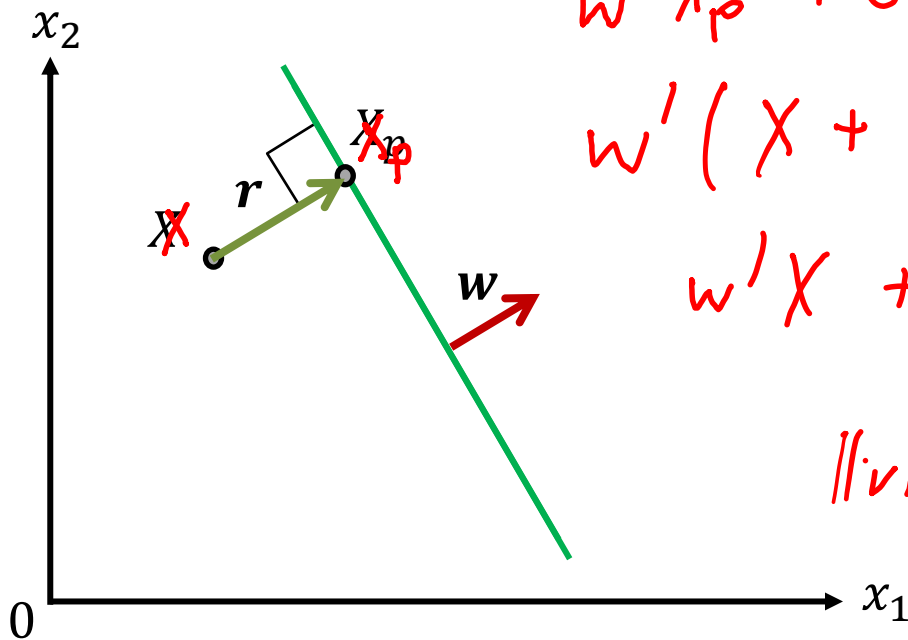
$$\mathbf{x}_p = \mathbf{x} + \|\mathbf{r}\| \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$\mathbf{w}'\mathbf{x}_p + b = 0$$

$$\mathbf{w}'\left(\mathbf{x} + \|\mathbf{r}\| \frac{\mathbf{w}}{\|\mathbf{w}\|}\right) + b = 0$$

$$\mathbf{w}'\mathbf{x} + \|\mathbf{r}\| \frac{\mathbf{w}'\mathbf{w} + b}{\|\mathbf{w}\|} = 0$$

$$\|\mathbf{r}\| = -\frac{(\mathbf{w}'\mathbf{x} + b)}{\|\mathbf{w}\|}$$



# Encoding the side using labels

- Training data is a collection  $\{\mathbf{x}_i, y_i\}$ ,  $i = 1, \dots, n$ , where each  $\mathbf{x}_i$  is an  $m$ -dimensional instance and  $y_i$  is the corresponding binary label encoded as  $-1$  or  $1$
- Given a **perfect** separation boundary,  $y_i$  will encode the side of the boundary each  $\mathbf{x}_i$  is on
- Thus the distance from the  $i$ -th point to a perfect boundary can be encoded as

$$\|\mathbf{r}_i\| = \frac{y_i(\mathbf{w}'\mathbf{x}_i + b)}{\|\mathbf{w}\|}$$

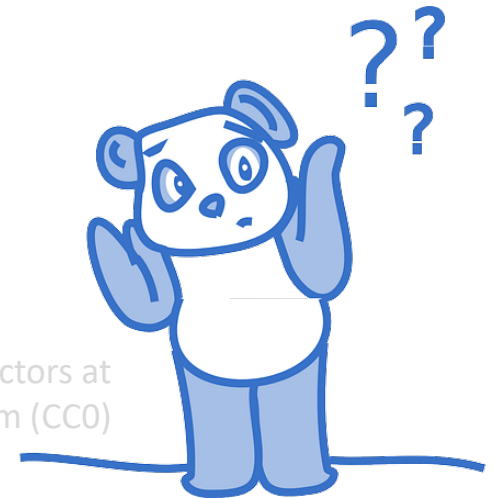


# Maximum margin objective

- The distance from the  $i$ -th point to a perfect boundary can be encoded as  $\|\mathbf{r}_i\| = \frac{y_i(\mathbf{w}'\mathbf{x}_i + b)}{\|\mathbf{w}\|}$
- The margin width is the distance to the closest point
- Thus SVMs aim to maximise  $\left( \min_{i=1,\dots,n} \frac{y_i(\mathbf{w}'\mathbf{x}_i + b)}{\|\mathbf{w}\|} \right)$  as a function of  $\mathbf{w}$  and  $b$

Do you see any problems with this objective?

art: OpenClipartVectors at  
pixabay.com (CC0)



# Constraining the objective

- SVMs aim to maximise  $\left( \min_{i=1,\dots,n} \frac{y_i(\mathbf{w}'\mathbf{x}_i+b)}{\|\mathbf{w}\|} \right)$
- Introduce (arbitrary) extra requirement  $\frac{y_{i^*}(\mathbf{w}'\mathbf{x}_{i^*}+b)}{\|\mathbf{w}\|} = \frac{1}{\|\mathbf{w}\|}$ 
  - \*  $i^*$  denotes index of a closest example to boundary
- SVM aims to find

$$\underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|$$

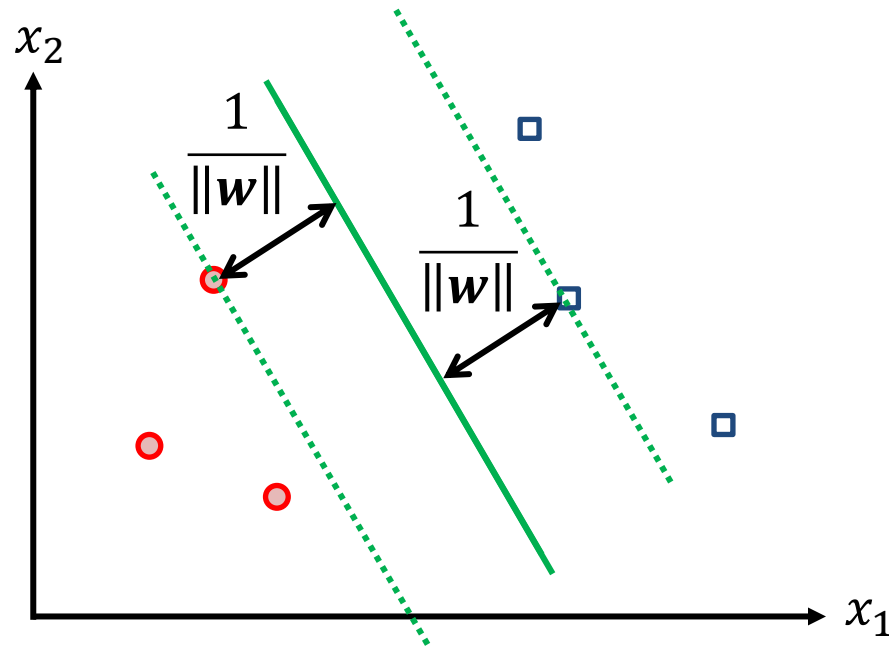
$$\text{s.t. } y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, n$$

# Hard margin SVM objective

We now have a major result: SVMs aim to find

$$\underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{w}\|$$

$$\text{s.t. } y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, n$$

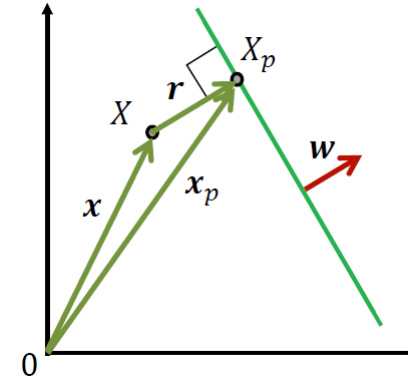
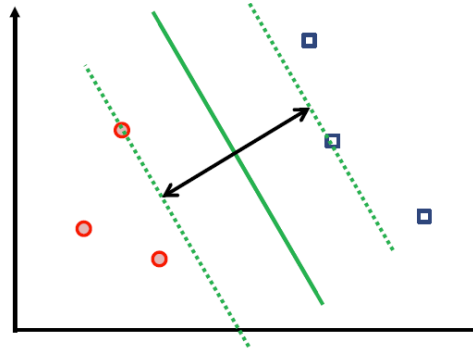
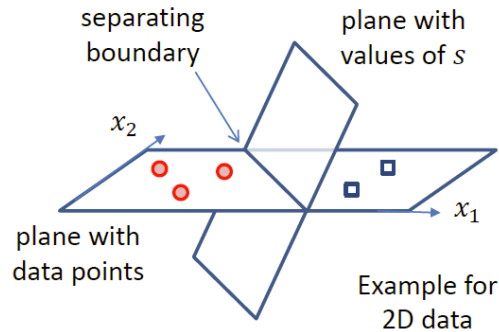


Note 1: parameter  $b$  is optimised indirectly by influencing constraints

Note 2: all points are enforced to be on or outside the margin

Therefore, this version of SVM is called *hard-margin SVM*

# Recap: hard-margin SVM



- SVM is a linear binary classifier
- Max margin: aim for boundary robust to noise
- Trick to resolve ambiguity  $\frac{y_i^*(w'x_i^* + b)}{\|w\|} = \frac{1}{\|w\|}$
- Hard-margin program:

$$\operatorname{argmin}_{w,b} \|w\| \text{ s.t. } y_i(w'x_i + b) \geq 1 \text{ for } i = 1, \dots, n$$

# SVM Objective as Regularised Loss

Relating the resulting objective function to  
that of other machine learning methods

# Previously in COMP90051 ...

1. Choose/design a model
2. Choose/design loss function
3. Find parameter values that minimise discrepancy on training data

How do SVMs fit this pattern?

# SVM as Regularised ERM

- Recall ridge regression objective

$$\text{minimise } \left( \sum_{i=1}^n (y_i - \mathbf{w}' \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|^2 \right)$$

- Hard margin SVM objective

data-dependent  
training error

$$\operatorname{argmin}_{\mathbf{w}} \|\mathbf{w}\|$$

data-independent  
regularisation term

$$\text{s.t. } y_i(\mathbf{w}' \mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, n$$

- The constraints can be interpreted as loss

$$l_{\infty} = \begin{cases} 0 & 1 - y_i(\mathbf{w}' \mathbf{x}_i + b) \leq 0 \\ \infty & 1 - y_i(\mathbf{w}' \mathbf{x}_i + b) > 0 \end{cases}$$

# Hard margin SVM loss

- The constraints can be interpreted as loss

$$l_{\infty} = \begin{cases} 0 & 1 - y_i(\mathbf{w}'\mathbf{x}_i + b) \leq 0 \\ \infty & 1 - y_i(\mathbf{w}'\mathbf{x}_i + b) > 0 \end{cases}$$

- In other words, for each point:
  - \* If it's on the right side of the boundary and at least  $\frac{1}{\|\mathbf{w}\|}$  units away from the boundary, we're OK, the loss is 0
  - \* If the point is on the wrong side, or too close to the boundary, we immediately give infinite loss thus prohibiting such a solution altogether

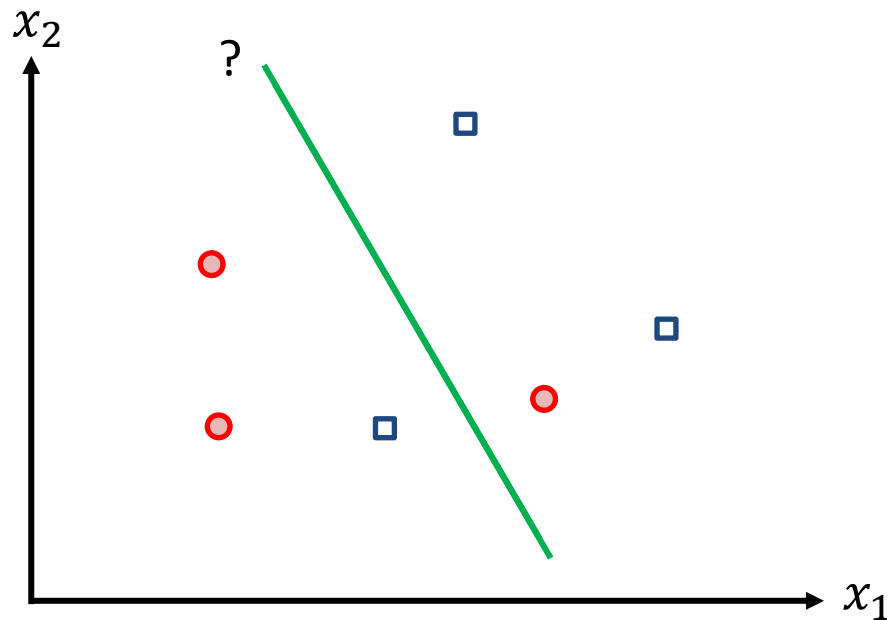


# Soft-Margin SVMs

Addressing linear inseparability

# When data is not linearly separable

- Hard-margin loss is too stringent (*hard!*)
- Real data is unlikely to be linearly separable
- If the data is not separable, hard-margin SVMs are in trouble

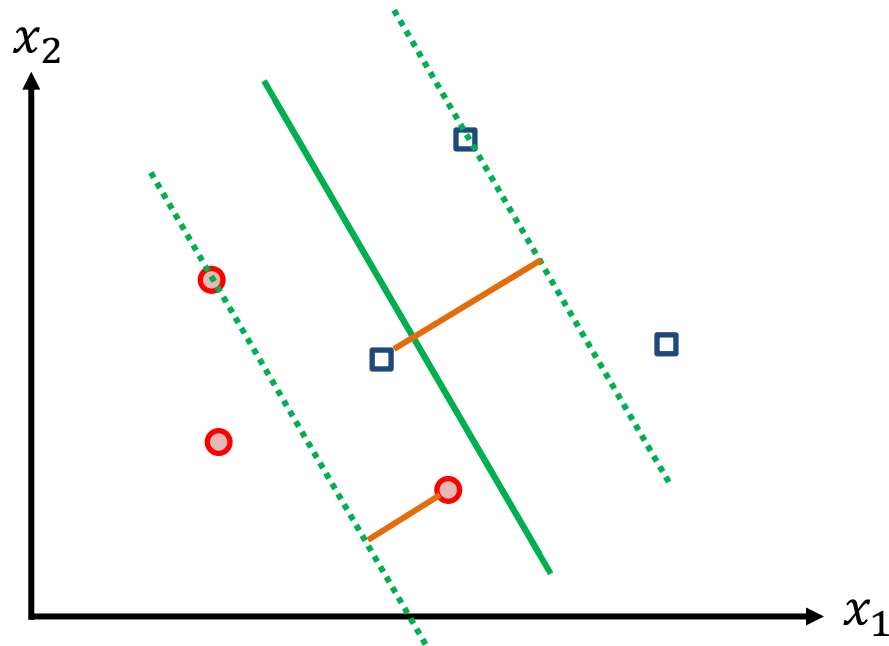


SVMs offer 3 approaches to address this problem:

1. *Still use hard-margin SVM, but **transform** the data (next lecture)*
2. ***Relax** the constraints (next slide)*
3. *The combination of 1 and 2 😊*

# Soft-margin SVM

- Relax constraints to allow points to be **inside the margin** or even on the **wrong side** of the boundary



However, we **penalise** boundaries by the extent of “violation”

In the figure, the objective penalty will take into account the orange distances

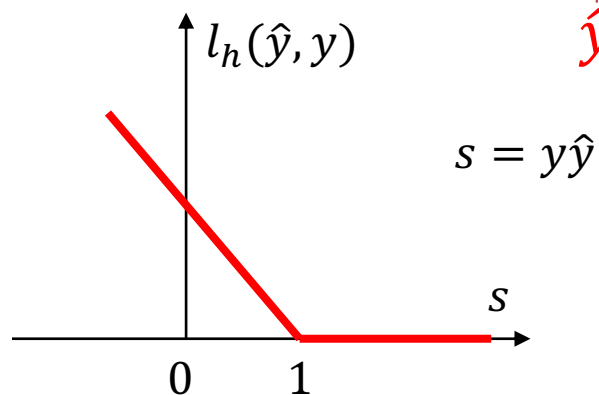
# Hinge loss: soft-margin SVM loss

- Hard-margin SVM loss

$$l_{\infty} = \begin{cases} 0 & 1 - y(\mathbf{w}'\mathbf{x} + b) \leq 0 \\ \infty & \text{otherwise} \end{cases}$$

- Soft-margin SVM loss (**hinge loss**)

$$l_h = \begin{cases} 0 & 1 - y(\mathbf{w}'\mathbf{x} + b) \leq 0 \\ 1 - y(\mathbf{w}'\mathbf{x} + b) & \text{otherwise} \end{cases}$$



compare this with  
perceptron loss

# Soft-margin SVM objective

- Soft-margin SVM **objective**

$$\operatorname{argmin}_{\mathbf{w}, b} \left( \sum_{i=1}^n l_h(\mathbf{x}_i, y_i, \mathbf{w}, b) + \lambda \|\mathbf{w}\|^2 \right)$$

- \* Reminiscent of ridge regression
- \* Hinge loss  $l_h = \max(0, 1 - y_i(\mathbf{w}'\mathbf{x}_i + b))$
- We are going to re-formulate this objective to make it more amenable to analysis

# Re-formulating soft-margin objective

- Define **slack variables** as an upper bound on loss

$$\xi_i \geq l_h = \max(0, 1 - y_i(\mathbf{w}'\mathbf{x}_i + b))$$

or equivalently  $\xi_i \geq 1 - y_i(\mathbf{w}'\mathbf{x}_i + b)$  and  $\xi_i \geq 0$

- Re-write the soft-margin SVM objective as:

$$\operatorname{argmin}_{\mathbf{w}, b, \xi} \left( \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \right)$$

$$\text{s.t. } \xi_i \geq 1 - y_i(\mathbf{w}'\mathbf{x}_i + b) \text{ for } i = 1, \dots, n$$

$$\xi_i \geq 0 \text{ for } i = 1, \dots, n$$

# Side-by-side: Two variations of SVM

- Hard-margin SVM objective\*:

$$\operatorname{argmin}_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{s.t. } y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, n$$

- Soft-margin SVM objective:

$$\operatorname{argmin}_{\mathbf{w}, b, \xi} \left( \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \right)$$

$$\text{s.t. } y_i(\mathbf{w}'\mathbf{x}_i + b) \geq 1 - \xi_i \text{ for } i = 1, \dots, n$$

$$\xi_i \geq 0 \text{ for } i = 1, \dots, n$$

- In the second case, the constraints are **relaxed (“softened”)** by allowing violations by  $\xi_i$ . Hence the name “soft margin”

\*Changed  $\|\mathbf{w}\|$  to  $0.5\|\mathbf{w}\|^2$  - monotonic increasing transform. Modified objective yields same solution.

# This lecture

- Support vector machines (SVMs) as maximum margin classifiers
- Deriving hard margin SVM objective
- SVM as regularised ERM
- Soft-margin SVM