

# **4F10: Deep Learning and Structured Data**

## **Support Vector Machines**

**José Miguel Hernández-Lobato**

Department of Engineering  
University of Cambridge

Michaelmas Term

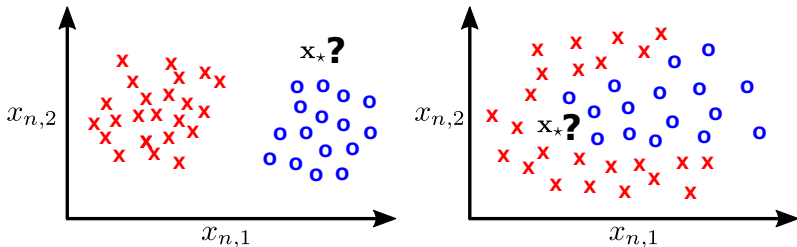
# The Problem of binary classification

Given a dataset  $\mathcal{D} = \{(x_n, t_n)\}_{n=1}^N$ , formed by pairs of input features  $x_n \in \mathbb{R}^d$  and target variables  $t_n \in \{-1, 1\}$ , we want to learn a classification function or **classifier**  $y(x)$  such that

$$\begin{aligned} y(x) &\geq 0 & \text{if } t_n = +1, \\ y(x) &< 0 & \text{if } t_n = -1. \end{aligned}$$

The classifier makes the correct prediction on a new input  $x_*$  when  $y(x_*)t_* > 0$ .

$$\circ \rightarrow t_n = 1 \quad \times \rightarrow t_n = -1$$



The **decision border** is the set of inputs for which  $y(x) = 0$ .

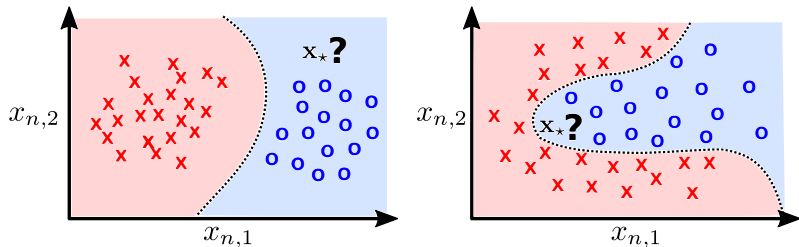
# The Problem of binary classification

Given a dataset  $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1}^N$ , formed by pairs of input features  $\mathbf{x}_n \in \mathbb{R}^d$  and target variables  $t_n \in \{-1, 1\}$ , we want to learn a classification function or **classifier**  $y(\mathbf{x})$  such that

$$\begin{aligned} y(\mathbf{x}) &\geq 0 & \text{if } t_n = +1, \\ y(\mathbf{x}) &< 0 & \text{if } t_n = -1. \end{aligned}$$

The classifier makes the correct prediction on a new input  $\mathbf{x}_\star$  when  $y(\mathbf{x}_\star)t_\star > 0$ .

$$\circ \rightarrow t_n = 1 \quad \times \rightarrow t_n = -1$$



The **decision border** is the set of inputs for which  $y(\mathbf{x}) = 0$ .

# Linear separability

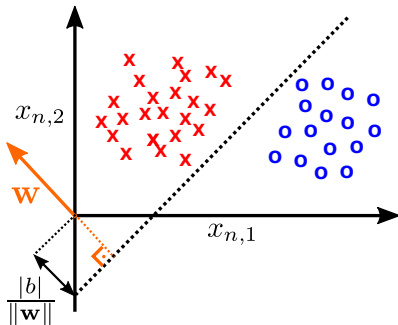
A classification problem is linearly separable when a classifier with a **linear decision border** makes no mistakes on the training data.

**Linear classifier:**

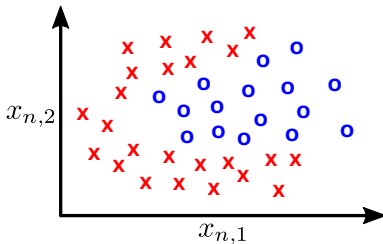
$$y(x) = \mathbf{w}^T \mathbf{x} + b,$$

$\mathbf{w}$  is a vector **orthogonal** to the decision border and  $b$  is the **bias** or intercept.

Linearly separable

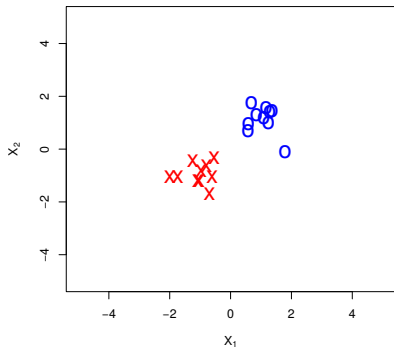


**Not** linearly separable



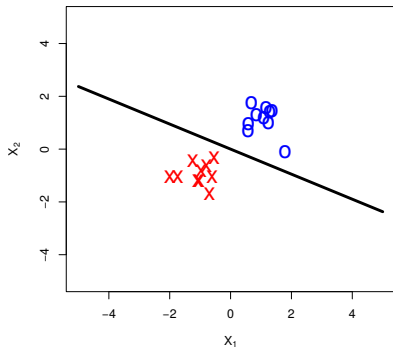
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



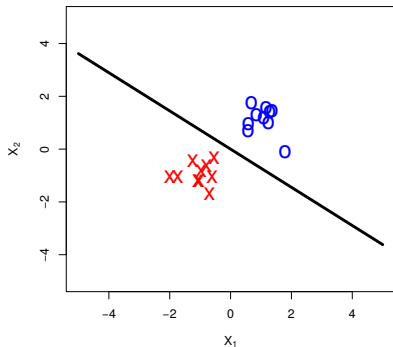
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



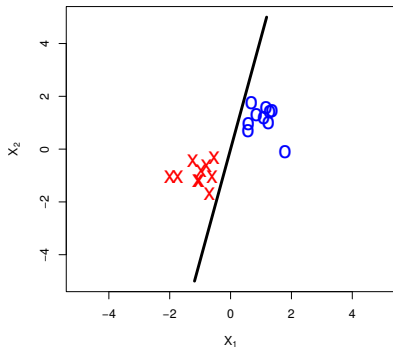
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



# Maximum margin classifiers

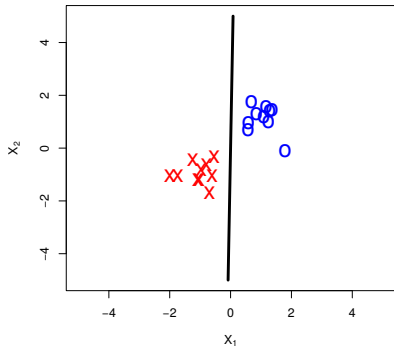
When the data is linearly separable, many possible  $w$  have zero training error.





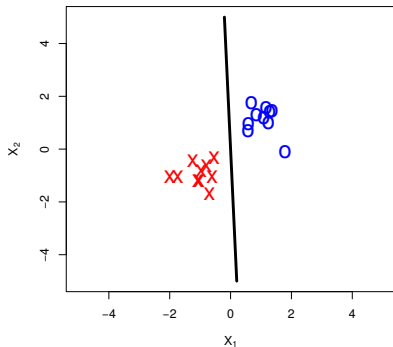
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



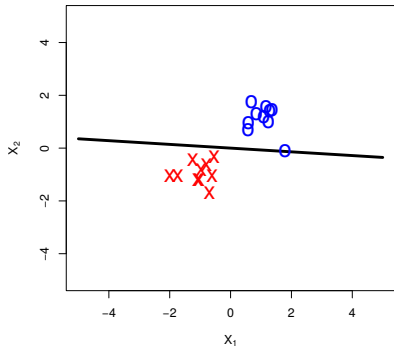
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



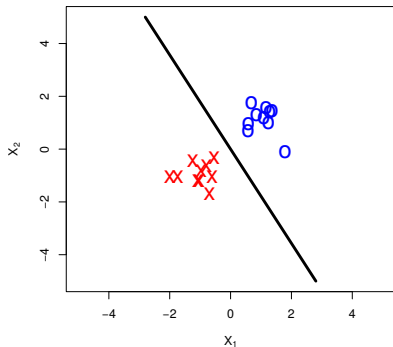
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



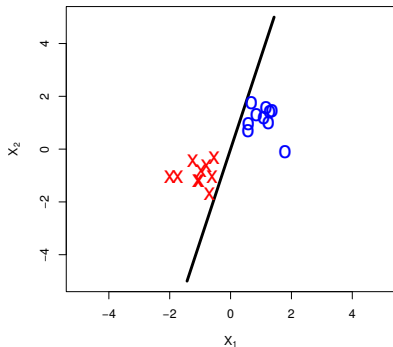
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



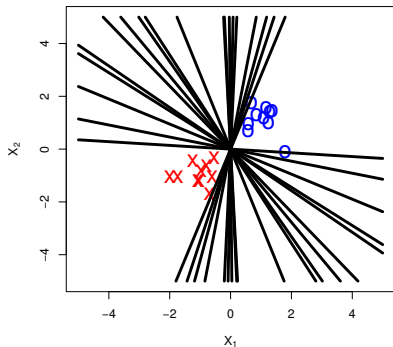
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



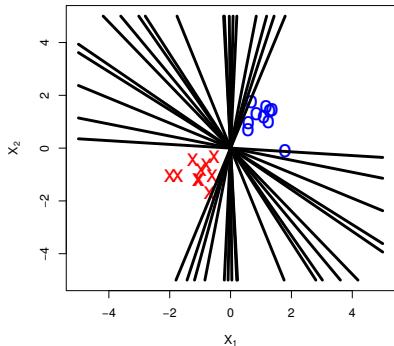
# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



# Maximum margin classifiers

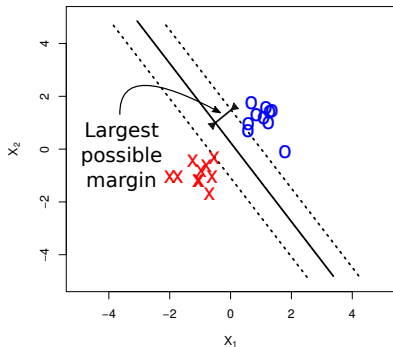
When the data is linearly separable, many possible  $w$  have zero training error.



**Maximum margin idea:** Choose plane whose distance to closest point in each class is maximal. This distance is called the **margin**.

# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.

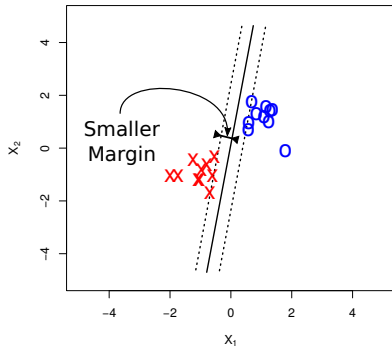
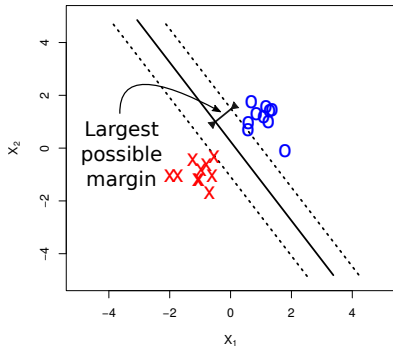


**Maximum margin idea:** Choose plane whose distance to closest point in each class is maximal. This distance is called the **margin**.



# Maximum margin classifiers

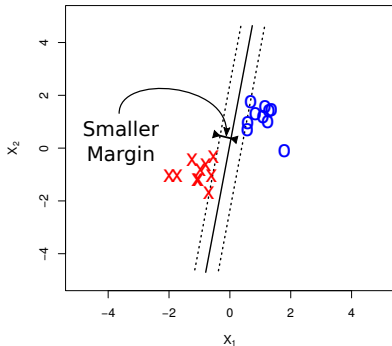
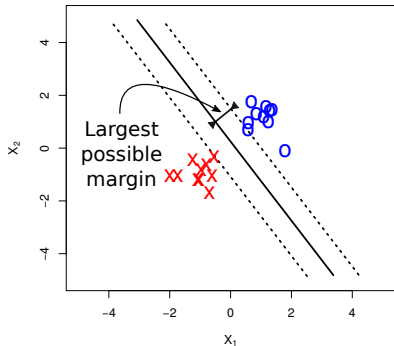
When the data is linearly separable, many possible  $w$  have zero training error.



**Maximum margin idea:** Choose plane whose distance to closest point in each class is maximal. This distance is called the **margin**.

# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



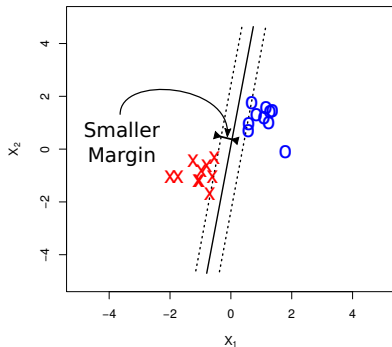
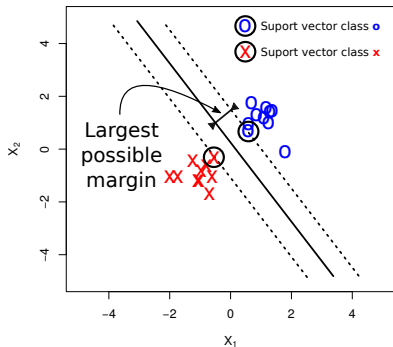
**Maximum margin idea:** Choose plane whose distance to closest point in each class is maximal. This distance is called the **margin**.

In a max-margin classifier, the data points on the margin are called **support vectors**. The max-margin classifier is fully determined by its support vectors.

Max margin-classifiers have some attractive theoretical justification.

# Maximum margin classifiers

When the data is linearly separable, many possible  $w$  have zero training error.



**Maximum margin idea:** Choose plane whose distance to closest point in each class is maximal. This distance is called the **margin**.

In a max-margin classifier, the data points on the margin are called **support vectors**. The max-margin classifier is fully determined by its support vectors.

Max margin-classifiers have some attractive theoretical justification.

# Computation of max-margin classifiers

The max-margin classifier or **support vector machine** (SVM) satisfies:

For any  $c > 0$ , any classifiers  $y(x) = w^T x + b$  and  $y'(x) = cy(x)$  have identical decision boundaries: the scale of the max-margin classifier is not important.

We can then choose the scale so that  $y(x_+) = +1$  for any positive support vector  $x_+$ . By symmetry,  $y(x_-) = -1$  for any negative support vector  $x_-$ . Therefore,

$$\begin{aligned}y(x_+) &= w^T x_+ + b = +1, \\y(x_-) &= w^T x_- + b = -1,\end{aligned}$$

The magnitude of the margin is then

$$\frac{w^T(x_+ - x_-)}{2\|w\|} = \frac{(w^T x_+ + b - w^T x_- - b)}{2\|w\|} = \frac{1}{\|w\|}. \quad (1)$$

At the solution, **(1) is maximal** and the data is correctly classified with the minimum value of  $|y(x)|$  achieved by support vectors. Therefore,  $t_n y(x_n) \geq 1 \ \forall n$ .

# Example solution

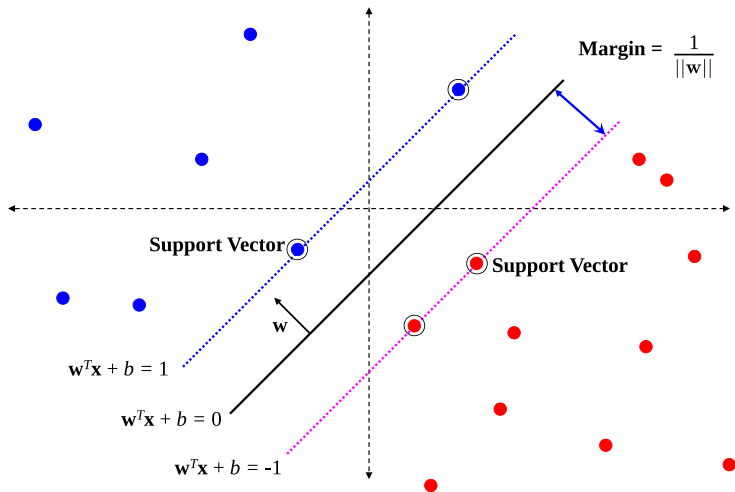


Figure source: A. Zisserman

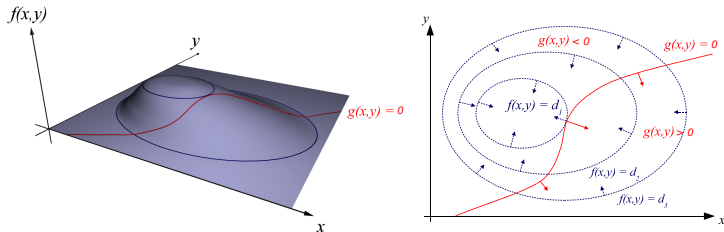
# Optimization with inequality constraints

Consider the constrained optimization problem

$$\text{maximize } f(x, y) \quad \text{subject to } g(x, y) \geq 0. \quad (2)$$

Let  $x_*$  and  $y_*$  be the solution to this problem.

Consider the constraint is instead  $g(x, y) = 0$ . We can then use Lagrange multipliers.



The solution must satisfy conditions  $\nabla_{x,y} f(x, y) = -\lambda \nabla_{x,y} g(x, y)$  and  $g(x, y) = 0$ .

Achieved by optimizing with respect to  $x$ ,  $y$  and  $\lambda$  the **Lagrangian** function:

$$\mathcal{L}(x, y, \lambda) = f(x, y) + \lambda g(x, y), \quad (3)$$

where  $\lambda \neq 0$  is called the **Lagrange multiplier**, which could be **positive or negative**.

# Karush-Kuhn-Tucker (KKT) conditions

When the constraint is  $g(x, y) \geq 0$ , we have two possibilities:

- ① **Constraint is not active at solution:**  $g(x_*, y_*) > 0$ .

Solution satisfies  $\nabla_{x,y} f(x, y) = 0$ . At the solution, the Lagrangian will also satisfy  $\nabla_{x,y} \mathcal{L}(x, y, \lambda) = 0$ , but only if  $\lambda = 0$ .

- ② **Constraint is active at solution:**  $g(x_*, y_*) = 0$ .

Stationary point of Lagrangian with  $\lambda \neq 0$ , that is,  $\nabla_{x,y} f(x, y) = -\lambda \nabla_{x,y} g(x, y)$ .

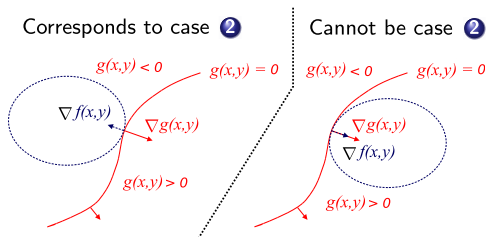
**$\lambda$  must be positive:**  $\lambda > 0$ , otherwise, solution located in region  $g(x, y) > 0$  and constraint would not be tight.

Solution to (2) obtained by solving  $\max_{x,y} \min_{\lambda} \mathcal{L}(x, y, \lambda)$  s.t.  $\lambda \geq 0$ .

In both cases above,  $\lambda g(x, y) = 0$  holds.

Solution to (2) satisfies **KKT conditions**:

$$g(x, y) \geq 0, \quad \lambda \geq 0, \quad \lambda g(x, y) = 0, \\ \nabla_{x,y} f(x, y) = -\lambda \nabla_{x,y} g(x, y).$$



# Optimization problem

Since maximizing  $1/\|w\|$  is the same as minimizing  $1/2\|w\|^2$ , the optimization problem can be written as

$$\min \quad \frac{1}{2}\|w\|^2 \quad \text{subject to} \quad t_n(w^T x_n + b) \geq 1, \quad n = 1, \dots, N.$$

This is a quadratic problem with linear constraints: **there is a unique minimum**.

By introducing **Lagrange multipliers**  $a_1 \geq 0, \dots, a_N \geq 0$  we obtain the objective

$$L(w, b, a) = \frac{1}{2}\|w\|^2 - \sum_{n=1}^N a_n \{t_n(w^T x_n + b) - 1\}. \quad (4)$$

The negative sign in  $a_n$  appears because we will be maximizing w.r.t. each  $a_n$ .

Setting gradients with respect to  $w$  and  $b$  to zero results in the following expressions:

$$w = \sum_{n=1}^N a_n t_n x_n, \quad 0 = \sum_{n=1}^N a_n t_n.$$

Using these expressions to eliminate  $w$  and  $b$  from (4) leads to the **dual problem**:

$$\max_a \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m x_n^T x_m \quad \text{s.t.} \quad \sum_{n=1}^N a_n t_n = 0, \quad \{a_n \geq 0\}_{n=1}^N.$$



# Solution and prediction formula

Need to optimize a **quadratic function** subject to **inequality constraints**.

**No analytic solution**, numerical methods have cost between  $\mathcal{O}(N^2)$  and  $\mathcal{O}(N^3)$ . SVMs are state-of-the-art methods for  $N \leq 50,000$ , otherwise too expensive!

Software for training SVMs: <http://www.support-vector-machines.org>.

From KKT conditions, at convergence, it holds that  $a_n(t_n y(x_n) - 1) = 0$ ,  $\forall n$ .

Hence, data points with  $a_n > 0$  satisfy  $t_n y(x_n) = 1$  and are **support vectors**.

After training, **predictions** can be made using

$$y(x) = \sum_{n \in \mathcal{S}} a_n t_n \mathbf{x}_n^T \mathbf{x} + b, \quad (5)$$

with  $\mathcal{S}$  being the set of indexes of support vectors:  $\mathcal{S} = \{n : a_n > 0\}$ .

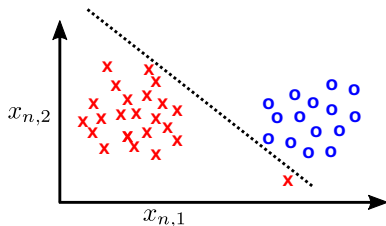
The bias  $b$  can be obtained from the set of constraints  $\{t_m y(x_m) = 1 : m \in \mathcal{S}\}$ :

$$b = \frac{1}{|\mathcal{S}|} \sum_{m \in \mathcal{S}} \left\{ t_m - \sum_{n \in \mathcal{S}} a_n t_n \mathbf{x}_n^T \mathbf{x}_m \right\}.$$

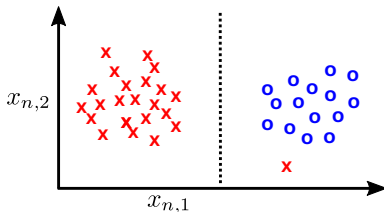
Numerical stability is improved by averaging the solution for  $b$  given by each constraint.

# Constraint violation and soft margin

Sometimes allowing for mistakes in the training data produces better classifiers.



The data points are linear separable, but the margin is very small.



We obtain a larger margin by ignoring the misclassified data point.

There is a **trade-off** between the margin size and the number of training errors.

Errors are allowed by introducing **slack variables**  $\xi_1, \dots, \xi_N$  in the constraints:

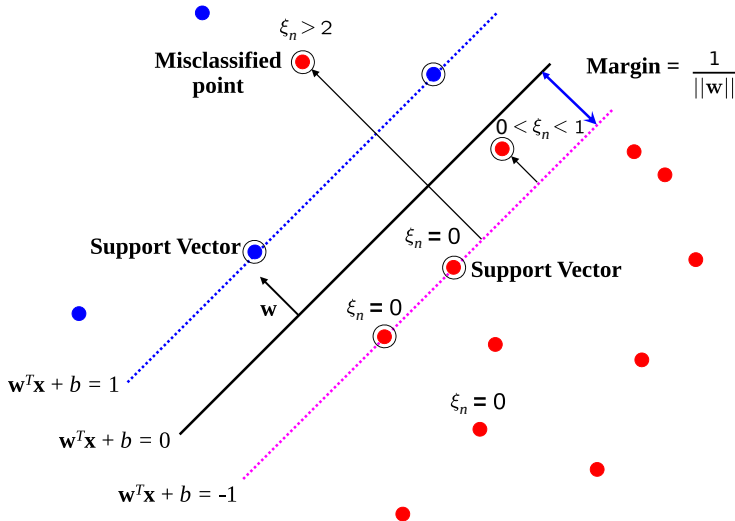
$$t_n(w^T x_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad n = 1, \dots, N,$$

Points with  $\xi_n = 0$  are correctly classified and are either on the margin or beyond.

Points with  $0 \leq \xi_n \leq 1$  lie inside the margin, but on the correct side.

Points with  $\xi_n > 1$  are misclassified.

# Example



# Optimization problem for soft margin

We now solve

$$\min \quad \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n \quad \text{s.t.} \quad t_n(w^T x_n + b) \geq 1 - \xi_n, \quad \xi_n \geq 0, \quad n = 1, \dots, N.$$

$C > 0$  controls the **trade-off** between the slack variable penalty and the margin.  
Small  $C$ : soft constraints, large margin. Large  $C$ : hard constraints, narrow margin.

By introducing **Lagrange multipliers**  $\{a_n \geq 0\}_{n=1}^N$ ,  $\{\mu_n \geq 0\}_{n=1}^N$  we obtain

$$L(w, b, a) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N a_n \{t_n(w^T x_n + b) - 1 + \xi_n\} - \sum_{n=1}^N \mu_n \xi_n.$$

Setting gradients with respect to  $w$ ,  $b$  and  $\{\xi_n\}$  to zero:  $w = \sum_{n=1}^N a_n t_n x_n$ ,  $0 = \sum_{n=1}^N a_n t_n$   $\{a_n = C - \mu_n\}_{n=1}^N$ .

Using this to eliminate  $w$ ,  $b$  and  $\{\xi_n\}$  from  $L(w, b, a)$  leads to the **dual problem**:

$$\max_a \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m x_n^T x_m \quad \text{s.t.} \quad \sum_{n=1}^N a_n t_n = 0, \quad \{0 \leq a_n \leq C\}_{n=1}^N.$$

# Prediction formula, support vectors and bias

Again, we solve a **quadratic function** subject to **inequality constraints**.

After training, **predictions** can be made using (5), as before.

At convergence, KKT conditions indicate that  $a_n(t_n y(x_n) - 1 + \xi_n) = 0, \forall n$ .

Hence, data points with  $a_n > 0$  satisfy  $t_n y(x_n) = 1 - \xi_n$  and are **support vectors**.

If  $0 < a_n < C$ , then  $\mu_n > 0$  and  $\xi_n = 0$  (from KKT conditions) and hence, such points lie on the margin.

Points with  $a_n = C$  can lie inside the margin boundaries and can either be

- correctly classified if  $\xi_n \leq 1$
- misclassified if  $\xi_n > 1$ .

$b$  obtained from constraints  $\{t_n y(x_n) = 1 - \xi_n : n \in \mathcal{S}\}$  with  $\xi_n = 0$  ( $0 < a_n < C$ ):

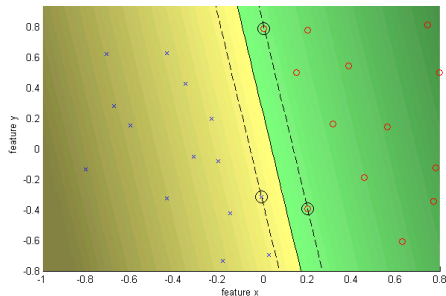
$$b = \frac{1}{|\mathcal{M}|} \sum_{n \in \mathcal{M}} \left\{ t_n - \sum_{m \in \mathcal{S}} a_m t_m x_m^T x_n \right\} \cdot \begin{array}{l} \text{We improve numerical stability by averaging} \\ \text{the solution for } b \text{ given by each constraint.} \end{array}$$

where  $\mathcal{M}$  contains the indexes of data points with  $0 < a_n < C$ .

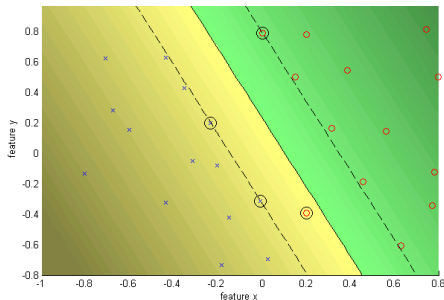
## Example

Data linearly separable but only with a narrow margin.

$C = \infty$ , hard margin.  
No training errors.



$C = 10$ , soft margin.  
One training error.



How to choose the optimal value for  $C$  in practice?

Solution: try many and choose one with optimal predictions on a validation set.

# Summary

When the data is linearly separable, many possible  $w$  have zero training error.

Max-margin classifier (MMC) has maximal distance to closest point in each class.

MMCs are theoretically justified and are expected to have low prediction error.

The MMC is found by solving a quadratic problem with inequality constraints.

Lagrange multipliers and the KKT conditions can be used to analyze the solution.

Slack variables can be used to obtain soft MMCs that allow for training errors.

Typically, there is a trade-off between  $\#$  of training errors and margin size.