

# Lecture 6: The SVM classifier

Slide credit: many of the slides were transcribed and adapted from **Andrew Zisserman** (“Machine Learning”, lecture 2) and **Yaser S. Abu-Mostafa** (“Learning from Data”, lecture 14)

# Outline

- Review of linear classifiers
  - Realizability == Linear separability
  - Revisiting perceptron
- Support Vector machine (SVM) classifier
  - Wide margin principle
  - Derivation
  - Dual form
  - Slack variables
  - Loss function
  - Multi class
  - Practical advice

# Binary Classification

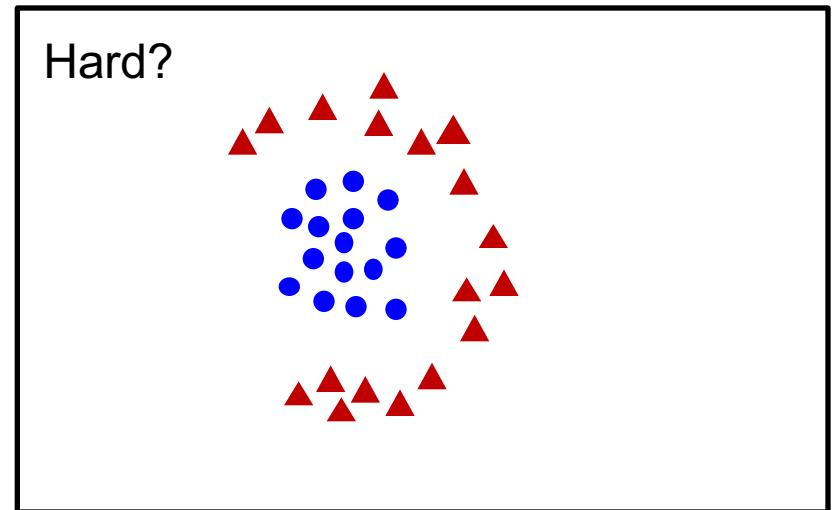
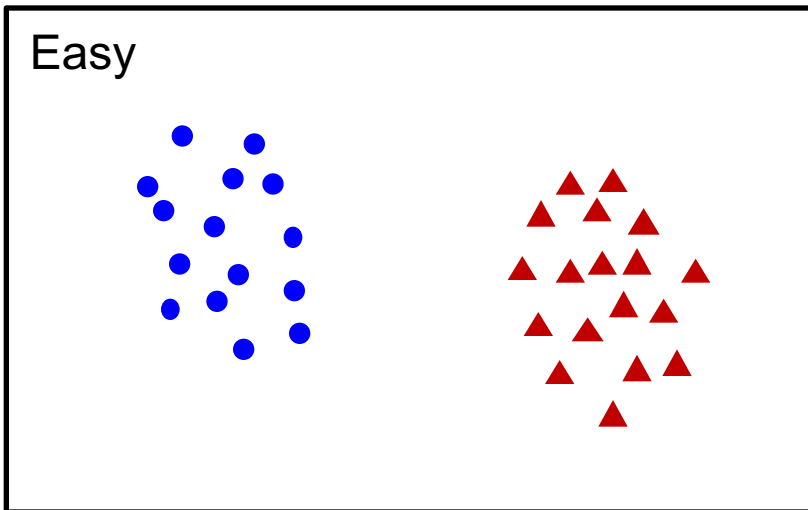
Given training data  $\{(x_i, y_i)\}_{i=1..N}$

with  $x_i \in \mathbb{R}^d$  and  $y_i \in \{-1, 1\}$

learn a classifier  $f(x)$  such that

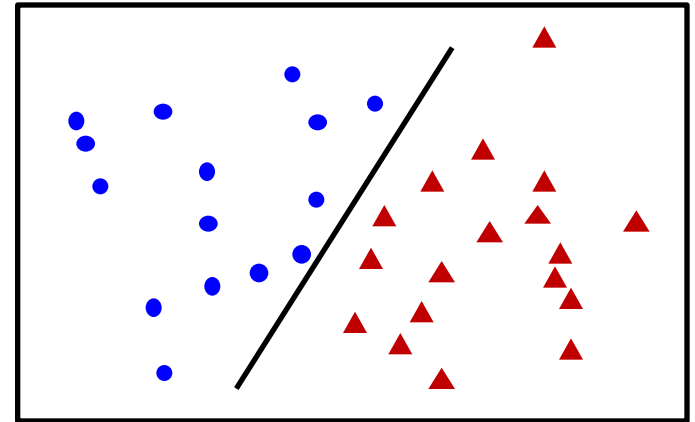
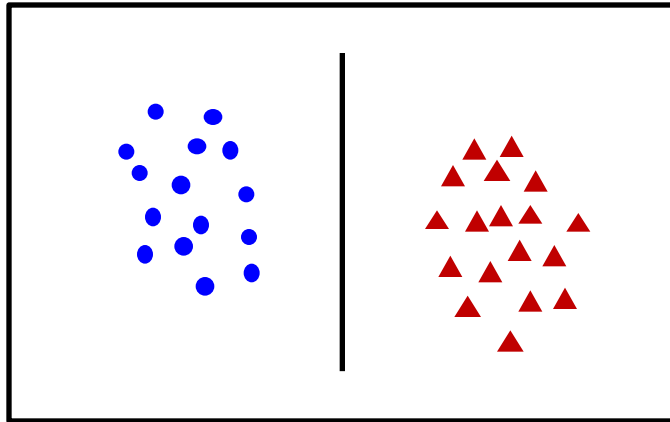
$$f(x_i) \begin{cases} \geq 0 & y_i = +1 \\ < 0 & y_i = -1 \end{cases}$$

i.e.,  $y_i f(x_i) > 0$  for a correct classification

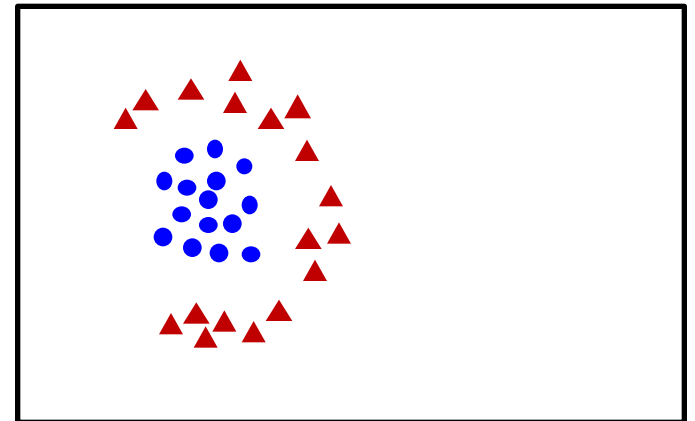
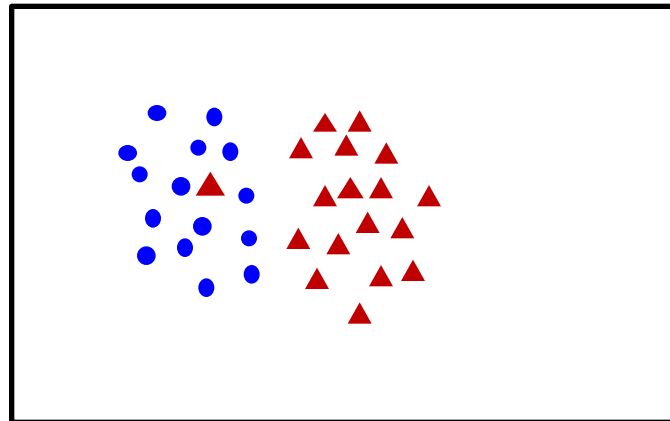


# Linear separability

linearly  
separable



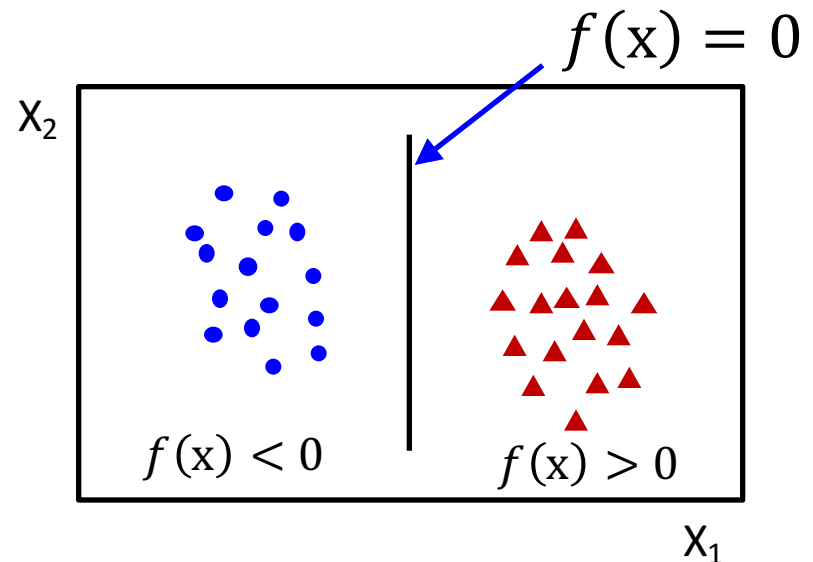
**not**  
linearly  
separable



# Linear classifiers

A linear classifier has the form

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

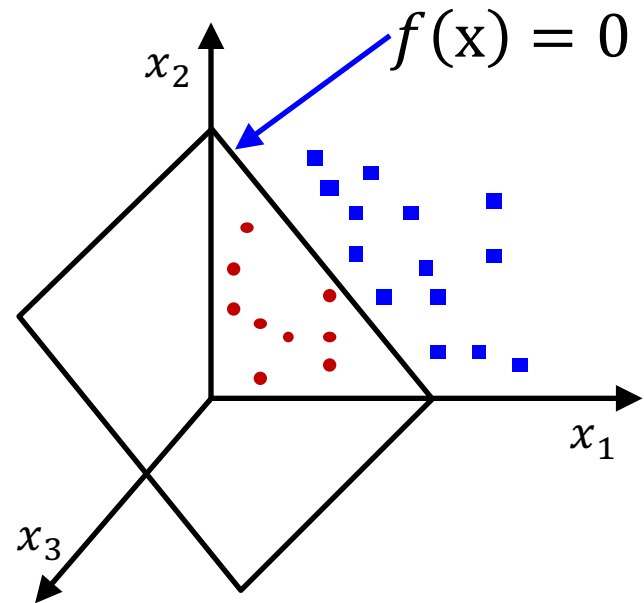


- in 2D the discriminant is a line
- $w$  is the **normal** to the line, and  $b$  the **bias**
- $w$  is known as the **weight vector**
- Before we added  $x_0=1$  and used  $w_0$  instead of  $b$

# Linear classifiers

A linear classifier has the form

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



- in 3D the discriminant is a plane, and in nD it is a hyperplane

# Reminder: The Perceptron Classifier

Given linearly separable data  $x_i$  labelled into two categories  $y_i = \{-1, 1\}$ , find a way to vector  $w$  such that the discriminant function

$$f(x_i) = w^T x_i + b$$

Separates the categories for  $i=1, \dots, N$

- How can we find this separating hyperplane?

## [The Perceptron Algorithm \(class 4\)](#)

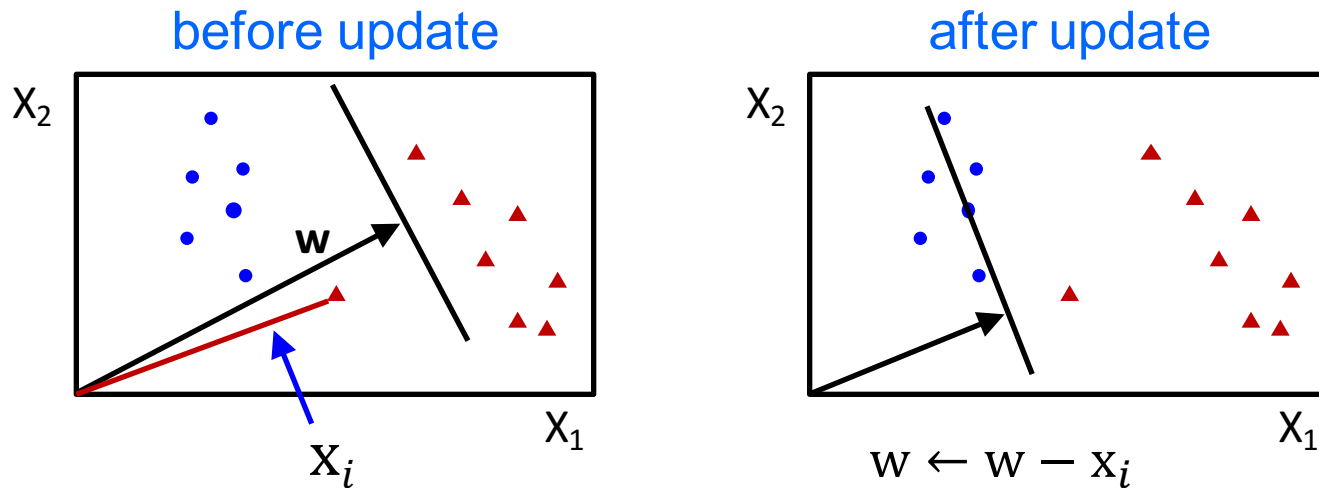
Write classifier as  $f(x_i) = \tilde{w}^T \tilde{x}_i + \omega_0 = w^T x_i$

where  $w = (\tilde{w}, \omega_0), x_i = (\tilde{x}_i, 1)$

- Initialize  $w = 0$
- Cycle through the data points  $\{x_i, y_i\}$ 
  - if  $x_i$  is misclassified then  $w \leftarrow w - \text{sign}(f(x_i))x_i$
- Until all the data is correctly classified

# For example in 2D

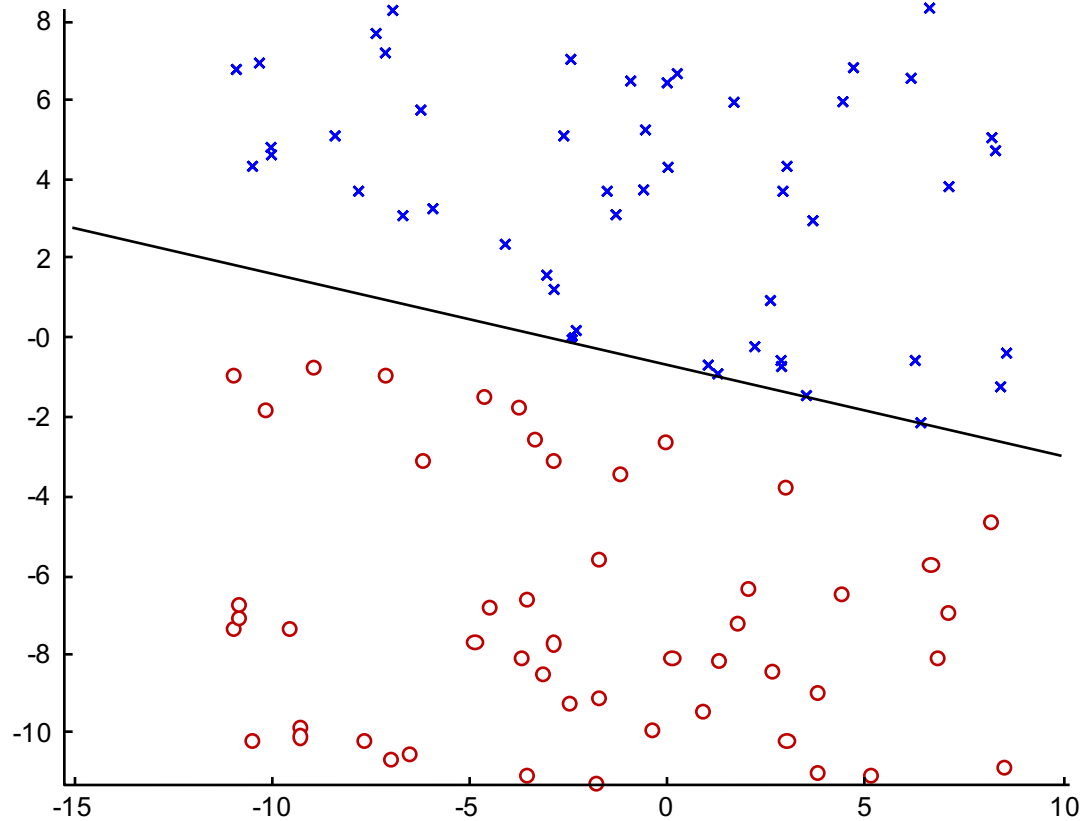
- Initialize  $w = 0$
- Cycle through the data points  $\{x_i, y_i\}$ 
  - If  $x_i$  is misclassified then  $w \leftarrow w + \text{sign}(f(x_i))x_i$
- Until all the data is correctly classified



After convergence  $w = \sum_i^N \alpha_i x_i$

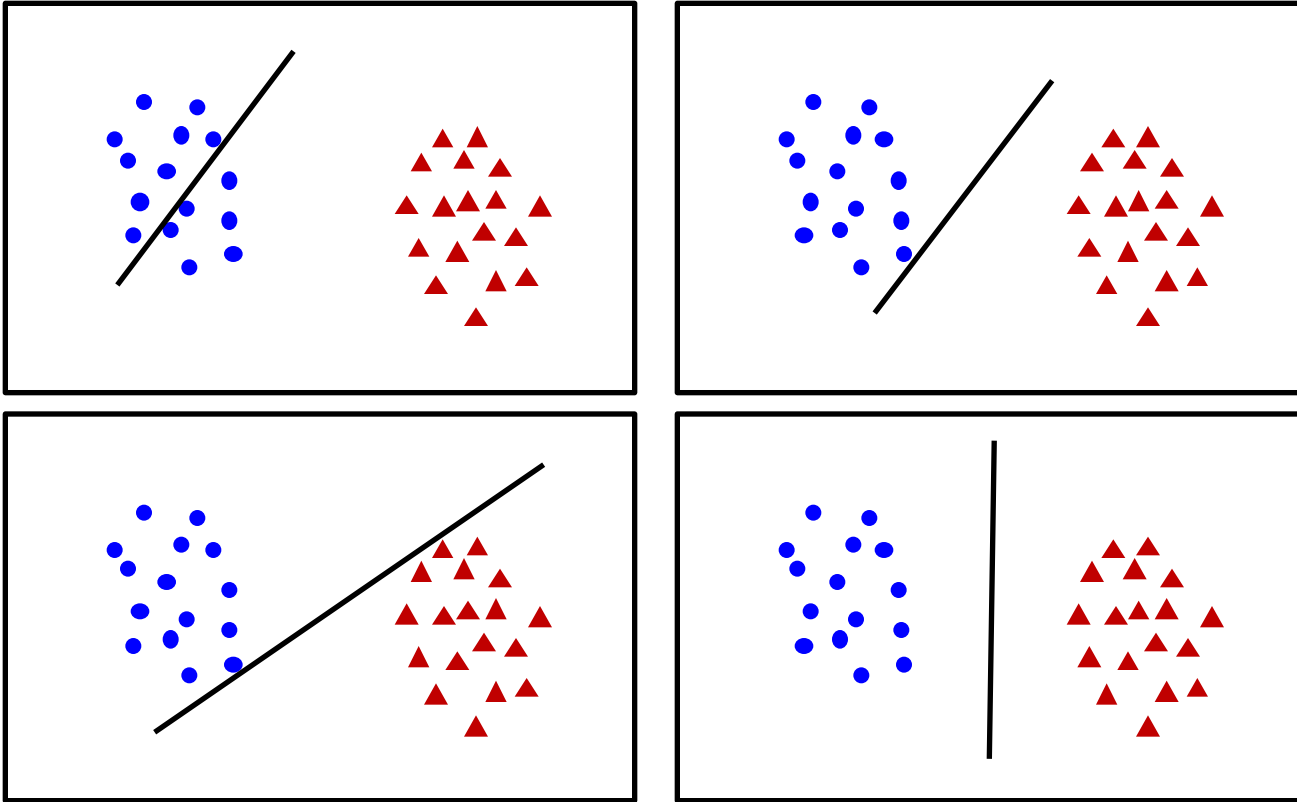


# Properties of the Perceptron Algorithm



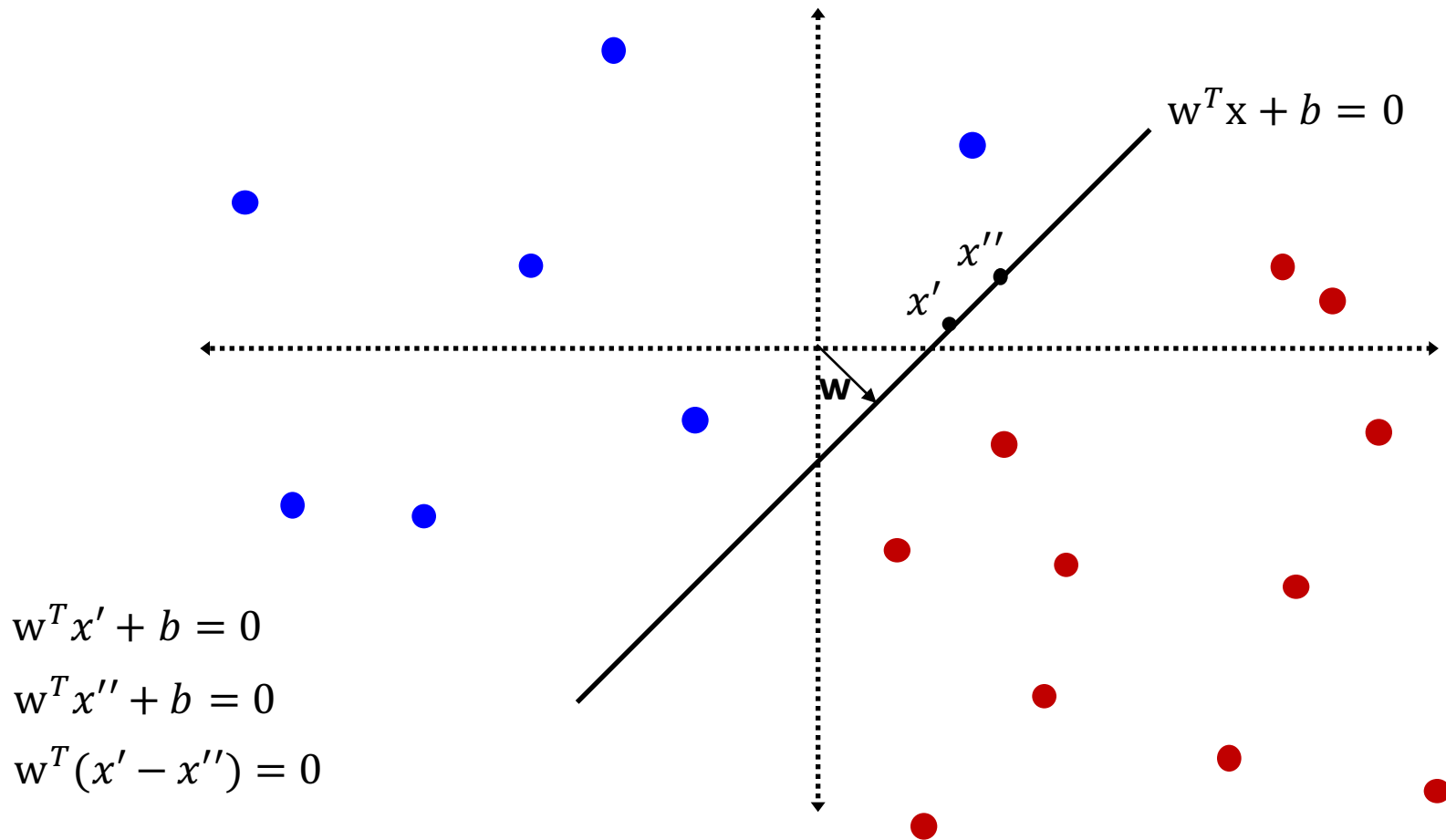
- Guaranteed convergence in the realizable case
- convergence can be very slow
- separating line close to training data

# What is the best $w$ ?

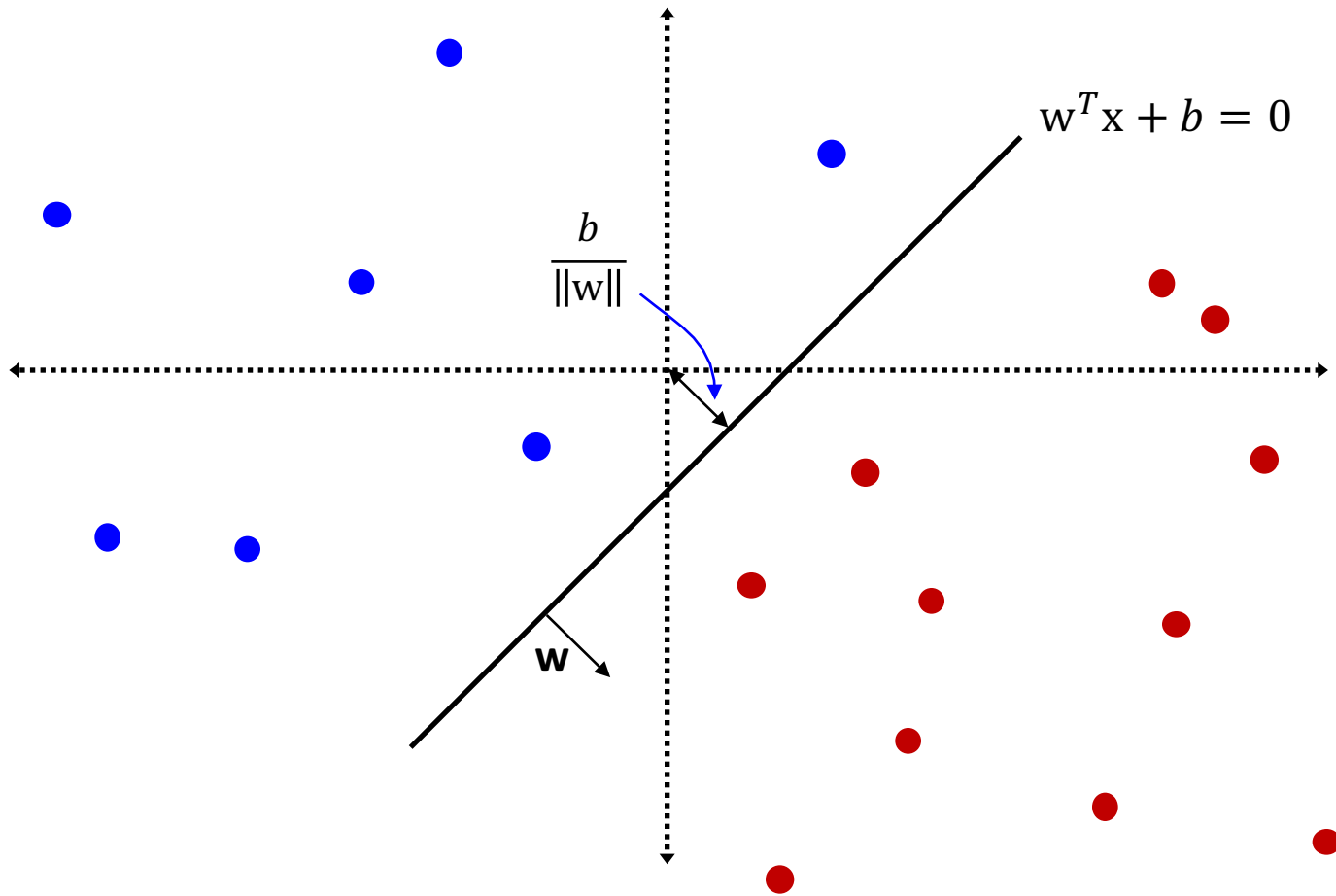


- **maximum margin** solution
  - most stable under perturbations of the inputs
  - more confident classification for all examples

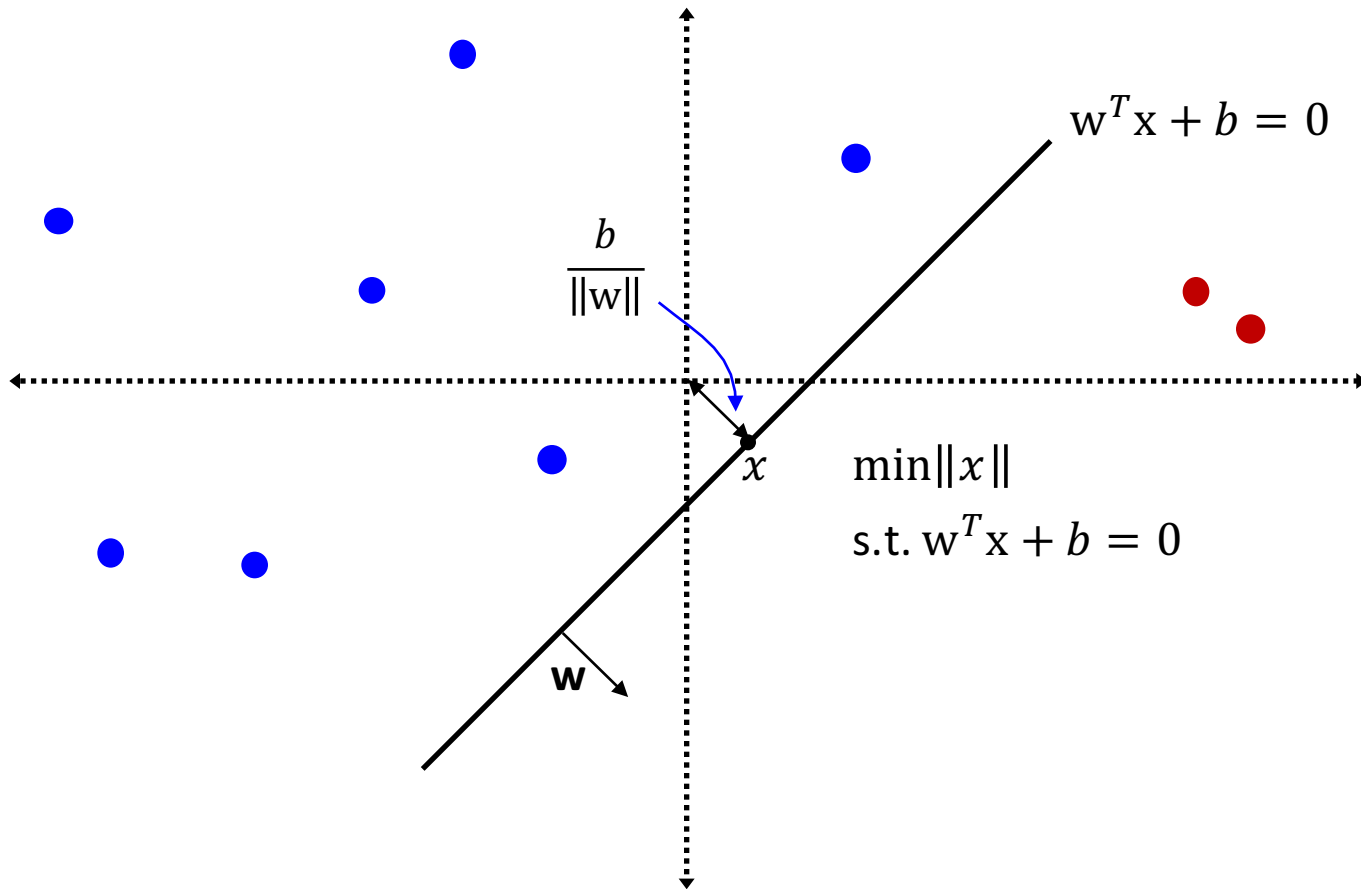
# What is $w$ ?



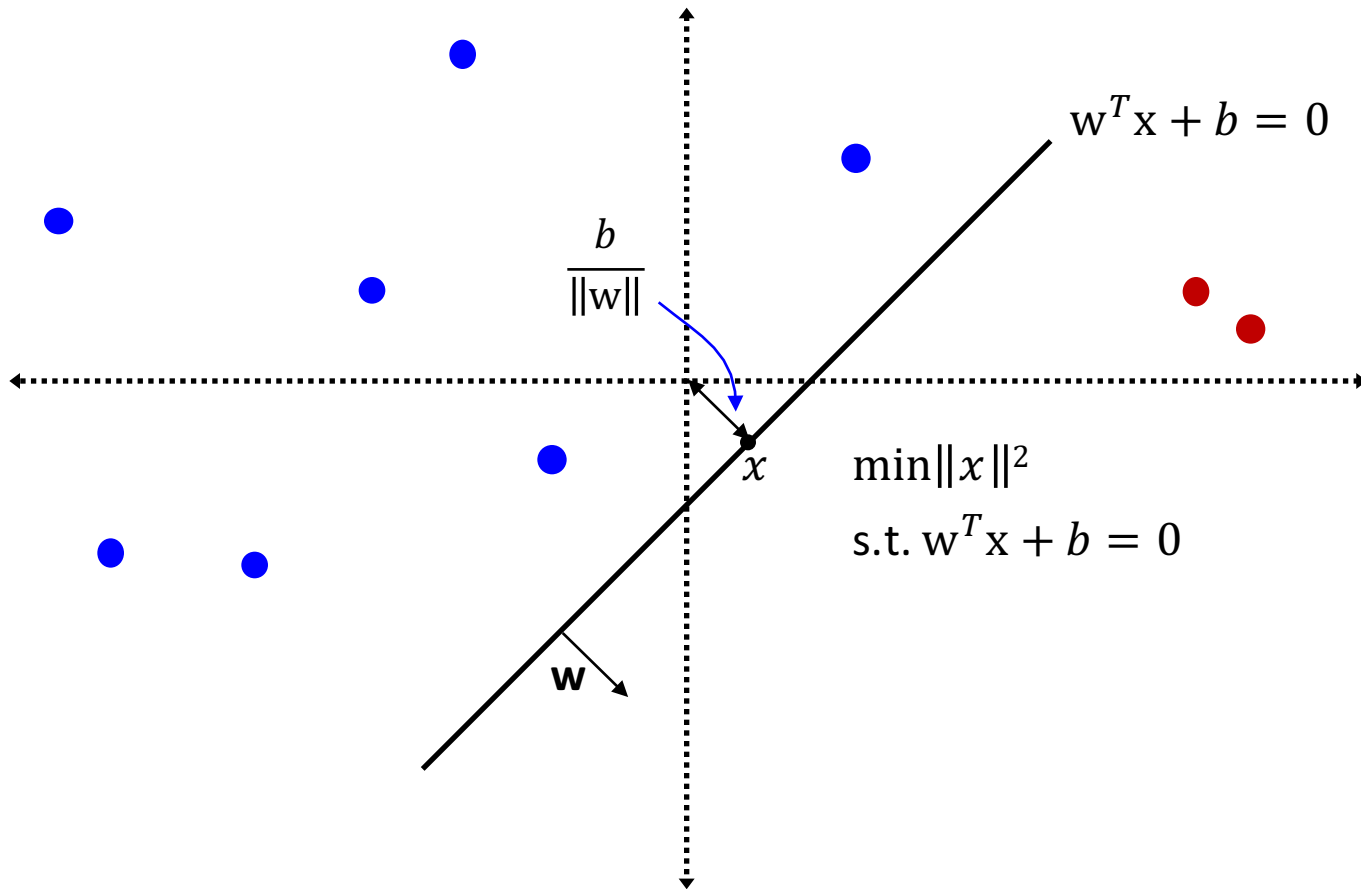
# What is $b$ ?



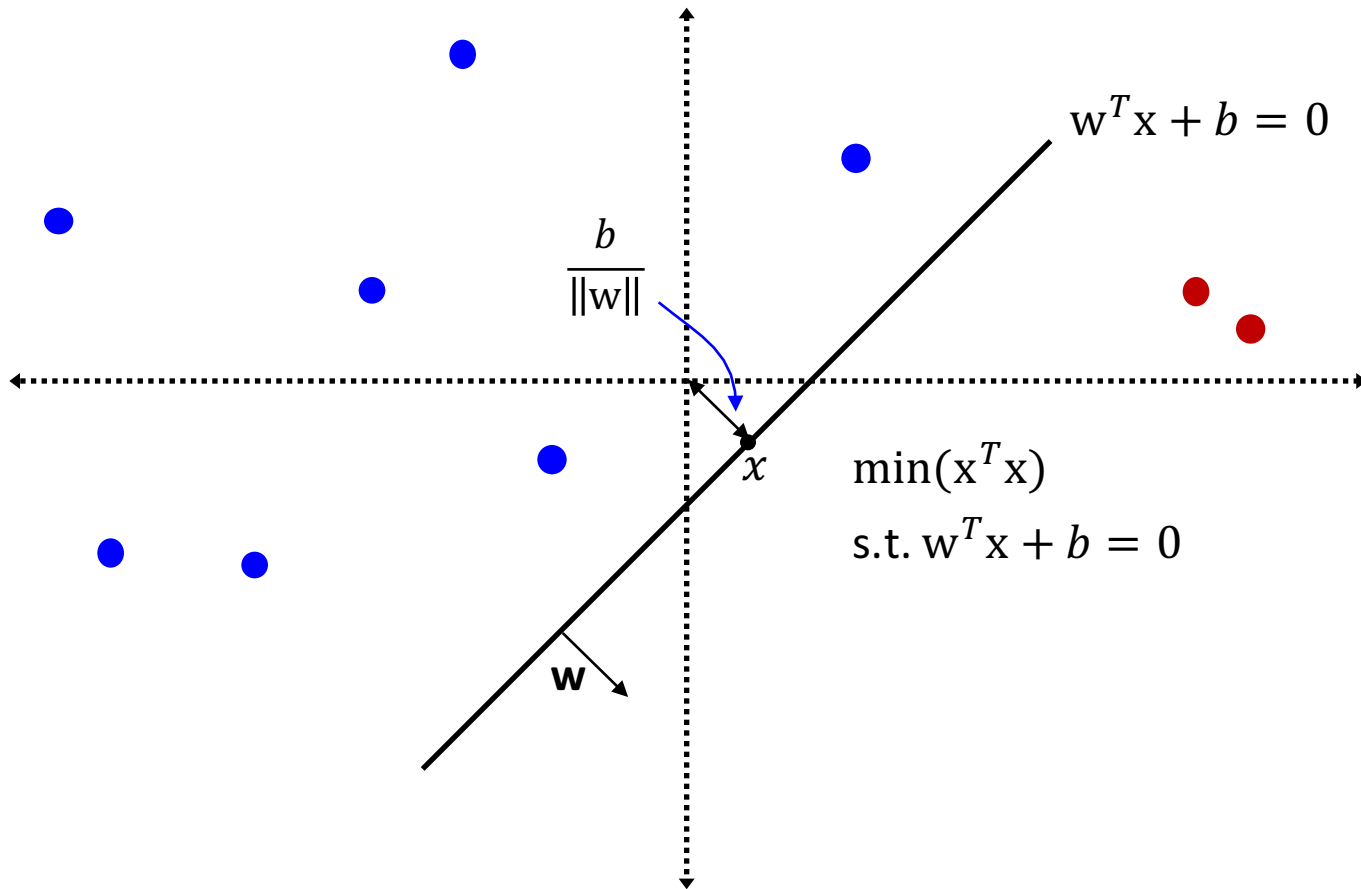
# What is $b$ ?



# What is $b$ ?



# What is $b$ ?



# What is b?

$$\min(x^T x)$$

$$\text{s.t. } w^T x + b = 0$$

$$L(x, \lambda) = x^T x - \lambda (w^T x + b)$$

$$2x - \lambda w = 0$$

$$x = \frac{\lambda}{2} w$$

$$w^T \left(\frac{\lambda}{2} w\right) + b = 0$$


$$\frac{\lambda}{2} w^T w + b = 0$$

$$\lambda = \frac{-2b}{w^T w}$$

$$x = \frac{\lambda}{2} w = \frac{-b}{w^T w} w$$

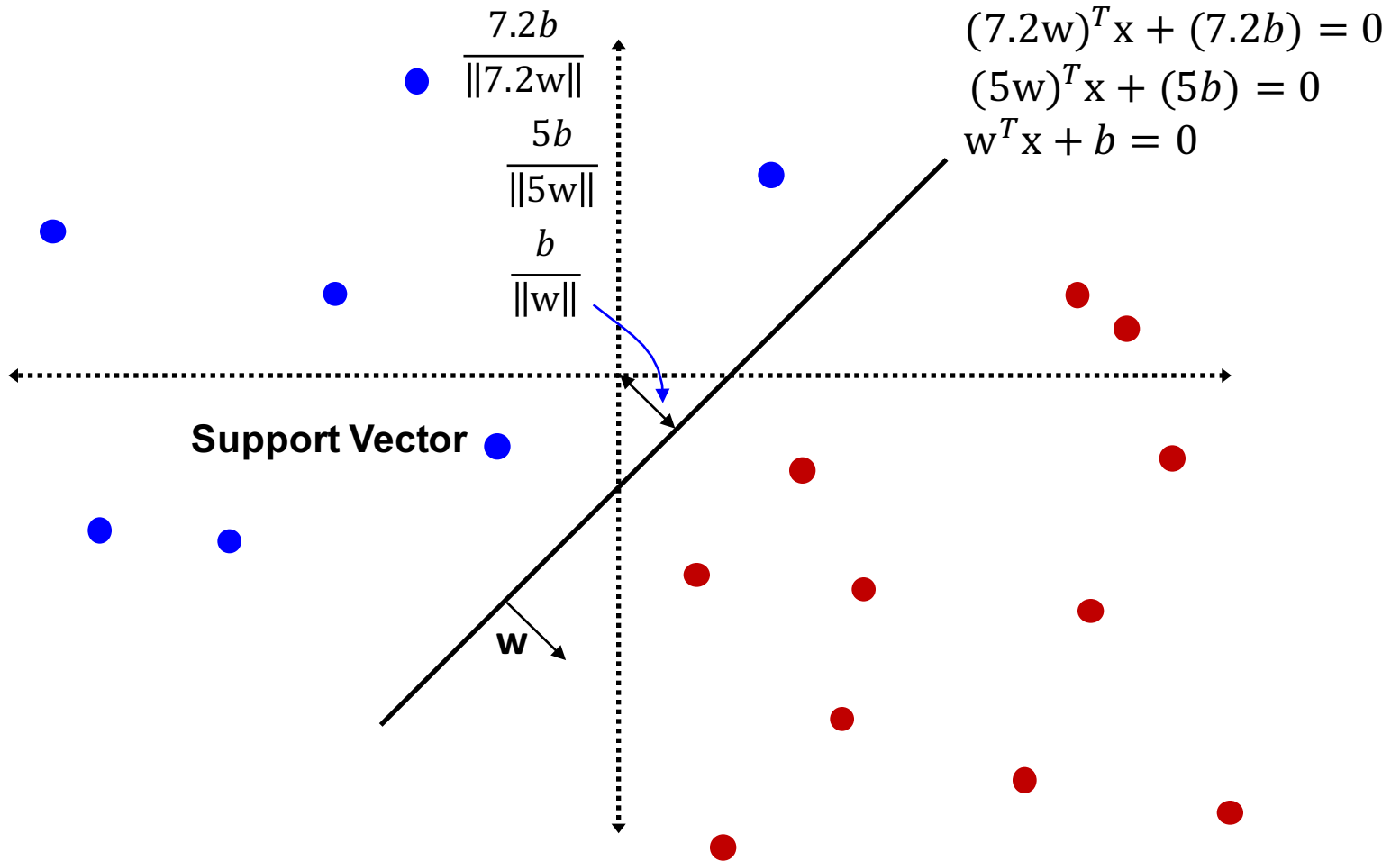
$$\|x\| = \sqrt{x^T x} = \sqrt{\left(\frac{-b}{w^T w} w\right)^T \left(\frac{-b}{w^T w} w\right)} = \frac{b}{w^T w} \sqrt{(w^T w)} = \frac{b}{\sqrt{(w^T w)}} = \frac{b}{\|w\|}$$

Google:  
"The matrix  
cookbook"





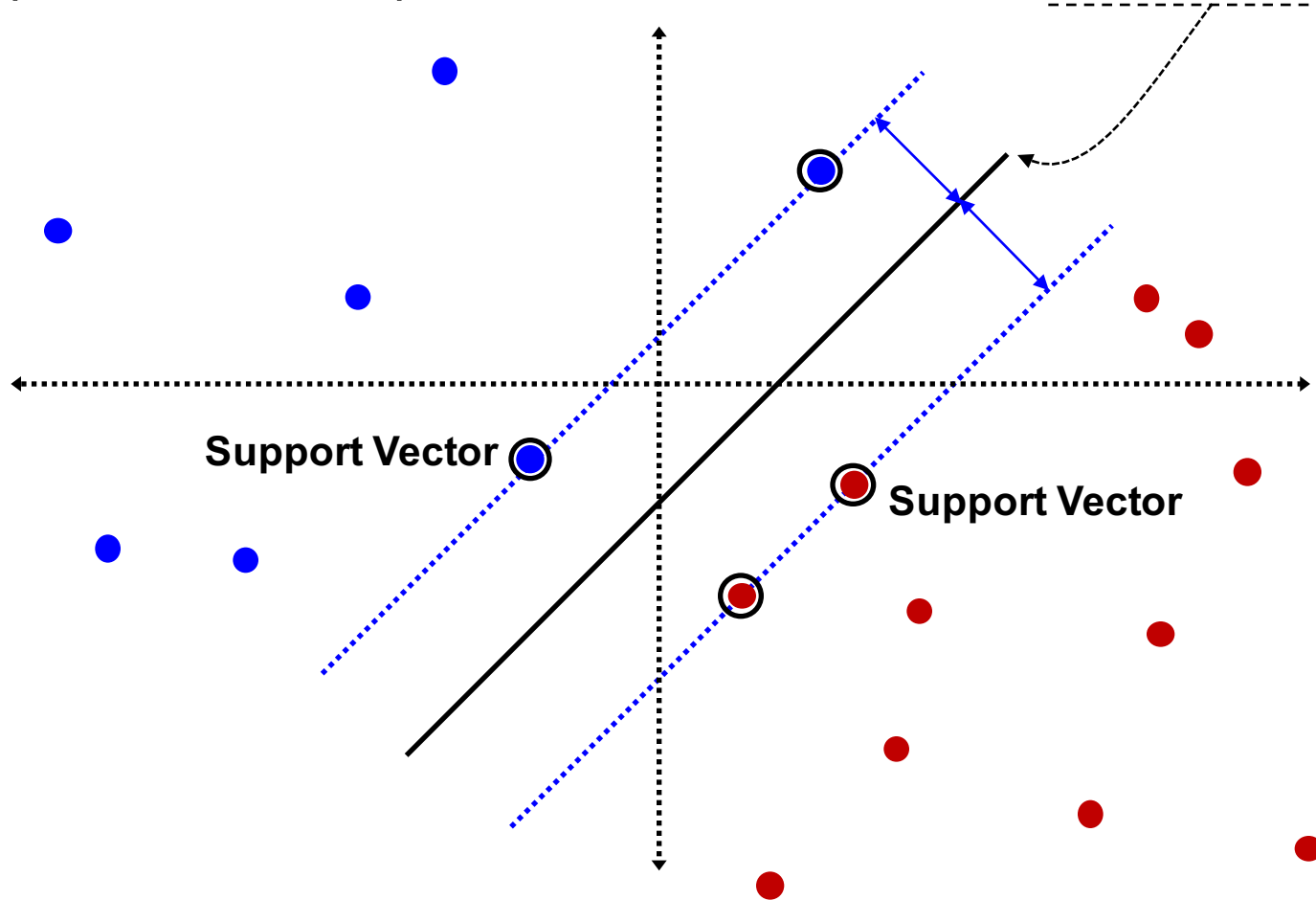
# Scale is arbitrary



What about negative scale?

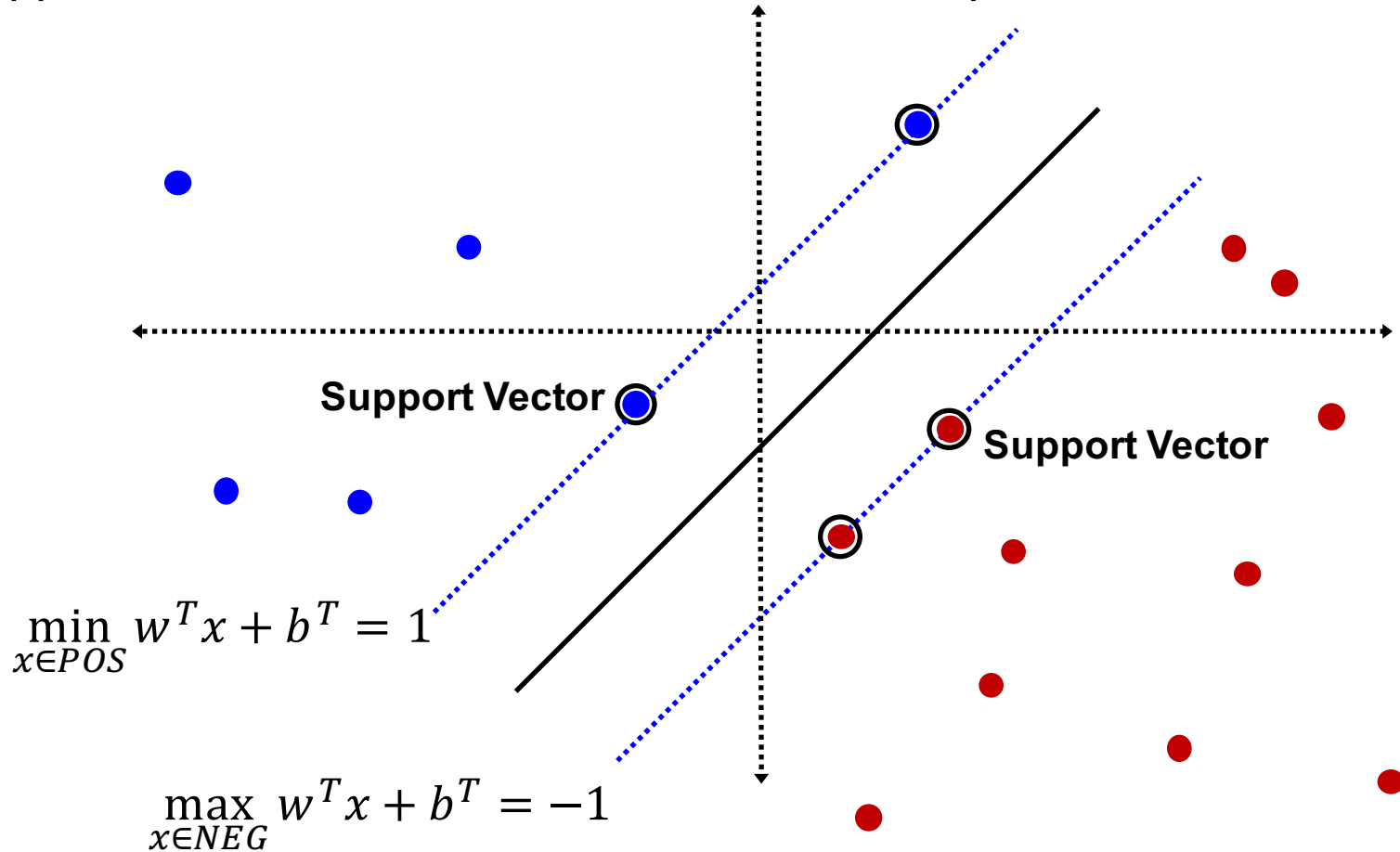
# Which points influence the line?

The *support vectors*, the points that are closest to the *decision boundary*



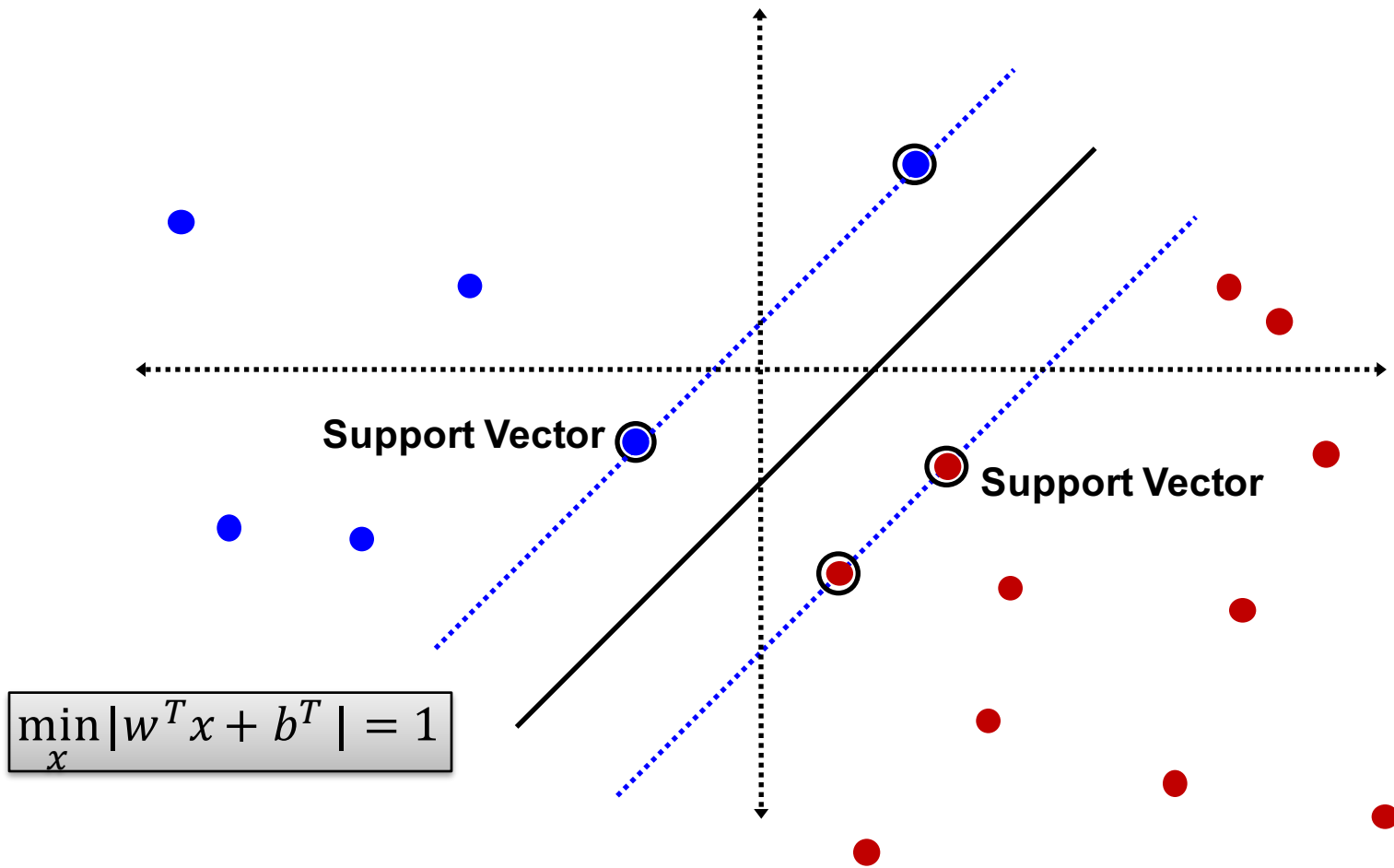
# Fixing the scale

Since we have freedom to choose the scale, we choose it such that support vectors have a value of 1 in the line equation



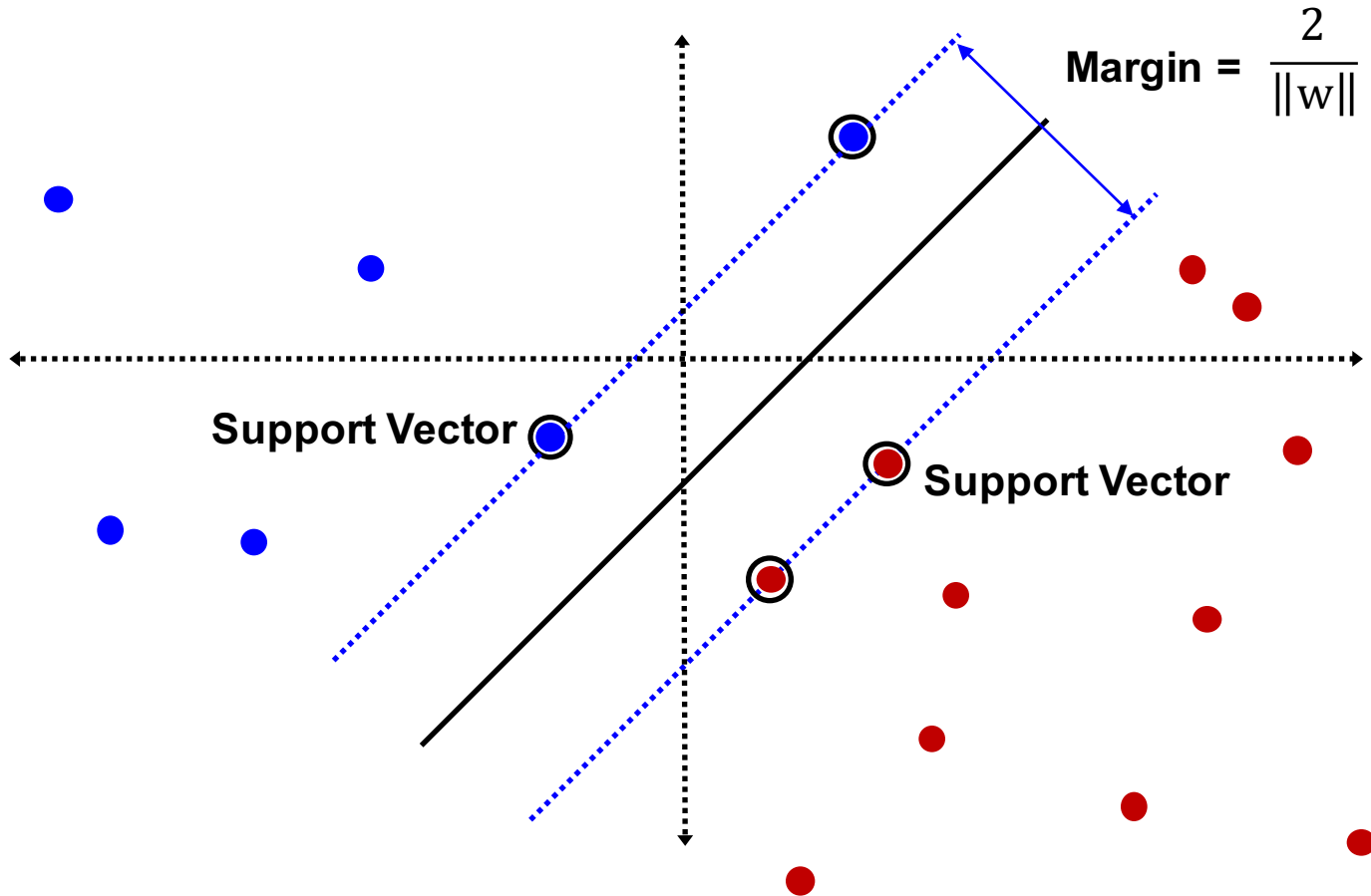
# Fixing the scale

Equivalently – all points have values  $\geq 1$



# Fixing the scale

In this case, we will show next, the margin is inversely proportional to  $\|w\|$



# Computing the distance

Let  $x_n$  be the nearest data point to the plane  $(w, b)$ . How far is it?

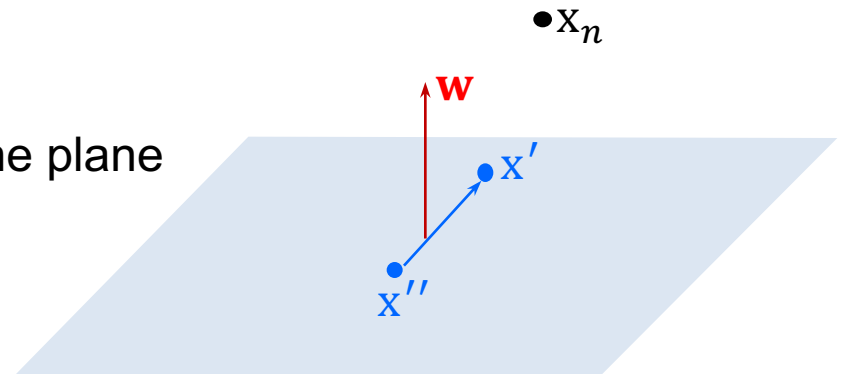
The distance between  $x_n$  and the plane  $w^T x + b = 0$   
where  $|w^T x_n + b| = 1$

The vector  $w$  is  $\perp$  to the plane:

As we've seen, take  $x'$  and  $x''$  on the plane

$$w^T x' + b = 0 \text{ and } w^T x'' + b = 0$$

$$\Rightarrow w^T (x' - x'') = 0$$



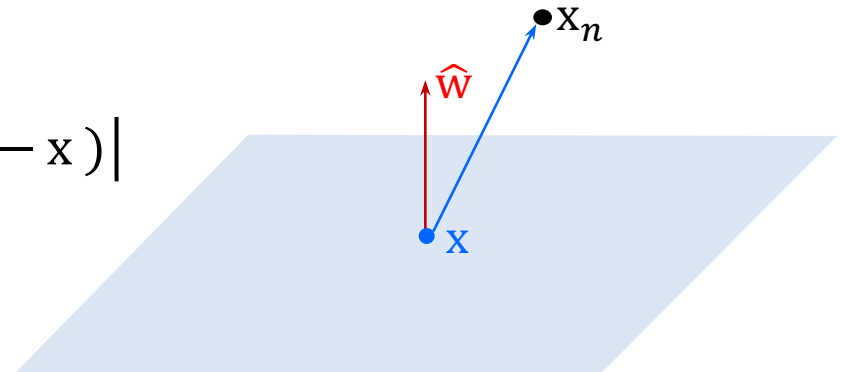
# and the distance is ...

Distance between  $x_n$  and the plane

Take any point  $x$  on the plane

Projection of  $x_n - x$  on  $w$

$$\hat{w} = \frac{w}{\|w\|} \Rightarrow \text{distance} = |\hat{w}^T (x_n - x)|$$



Distance =

$$\frac{1}{\|w\|} |w^T x_n - w^T x| = \frac{1}{\|w\|} | \underline{w^T x_n + b} - \underline{w^T x - b} |$$

# The vanilla SVM optimization problem

Maximize  $\frac{1}{\|w\|}$

subject to  $\min_{n=1,2,\dots,N} |w^T x_n + b| = 1$



# The vanilla SVM optimization problem

$$\text{Maximize } \frac{1}{\|\mathbf{w}\|}$$

$$\text{subject to } \min_{n=1,2,\dots,N} |\mathbf{w}^T \mathbf{x}_n + b| = 1$$

$$\text{Notice: } |\mathbf{w}^T \mathbf{x}_n + b| = y_n (\mathbf{w}^T \mathbf{x}_n + b)$$

$$\text{Maximize } \frac{1}{\|\mathbf{w}\|}$$

$$\text{subject to } y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for } n = 1, 2, \dots, N$$

$$\text{Minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{subject to } y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \text{ for } n = 1, 2, \dots, N$$

# Constrained optimization

Minimize  $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

subject to  $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$  for  $n = 1, 2, \dots, N$

$\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$

This is a quadratic optimization problem subject to linear constraints  
a.k.a Quadratic Programming

$\mathbf{x} = \text{quadprog}(\mathbf{H}, \mathbf{f}, \mathbf{A}, \mathbf{b})$  minimizes  $\frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{f}^T \mathbf{x}$

subject to the restrictions  $\mathbf{A} \mathbf{x} \leq \mathbf{b}$ .

A is a matrix of doubles, and b is a vector of doubles.

What is  $\mathbf{x}$ ? What is  $\mathbf{H}$ ? What if  $\mathbf{f}$ ?

H is positive semi-definite  $\rightarrow$  unique global solution

# Constrained optimization

Minimize  $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

subject to  $y_n (\mathbf{w}^T x_n + b) \geq 1$  for  $n = 1, 2, \dots, N$

$\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$

Lagrange?      inequality constraints

# Lagrange formulation

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

**Minimize** w.r.t  $\mathbf{w}$  and  $b$  and **Maximize** w.r.t each  $\alpha_n \geq 0$

$\alpha_n$  plays the adversary:

If for all  $n$ ,  $y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1 > 0$

then  $\max_{\boldsymbol{\alpha}} L(\mathbf{w}, b, \boldsymbol{\alpha})$  is obtained at  $\boldsymbol{\alpha} = \mathbf{0}$  and is simply  $\frac{1}{2} \mathbf{w}^T \mathbf{w}$

If for any  $n$ ,  $y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1 < 0$

then  $\alpha_n = \infty$  and  $\max_{\boldsymbol{\alpha}} L(\mathbf{w}, b, \boldsymbol{\alpha})$  becomes  $\infty$

# Lagrange formulation

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

**Minimize** w.r.t  $\mathbf{w}$  and  $b$  and **Maximize** w.r.t each  $\alpha_n \geq 0$

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0$$

$$\frac{\partial L}{\partial b} = - \sum_{n=1}^N \alpha_n y_n = 0$$

# Substituting...

$$\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n \quad \text{and employing} \quad \sum_{n=1}^N \alpha_n y_n = 0$$

in the Lagrangian

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

we get

$$L(\alpha) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \mathbf{x}_n^T \mathbf{x}_m + \sum_{n=1}^N \alpha_n$$

Maximize w.r.t to  $\alpha$  subject to  $\alpha_n \geq 0$  for  $n = 1, \dots, N$  and  $\sum_{n=1}^N \alpha_n y_n = 0$

# The solution – quadratic programming

$$\min_{\alpha} \frac{1}{2} \alpha^T \underbrace{\begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \dots & y_1 y_N x_1^T x_N \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 & \dots & y_2 y_N x_2^T x_N \\ \dots & \dots & \dots & \dots \\ y_N y_1 x_N^T x_1 & y_N y_2 x_N^T x_2 & \dots & y_N y_N x_N^T x_N \end{bmatrix}}_{\text{quadratic coefficients}} \alpha + \underbrace{(-1^T)}_{\text{linear}} \alpha$$

subject to  $\underbrace{y^T \alpha = 0}_{\text{linear constraint}}$

$$\underbrace{0}_{\text{lower bounds}} \leq \alpha \leq \underbrace{\infty}_{\text{upper bounds}}$$

Minimize  $\frac{1}{2} w^T w$

subject to  $y_n (w^T x_n + b) \geq 1$   
 $n = 1, 2, \dots, N$

$w \in \mathbb{R}^d, b \in \mathbb{R}$

# QP hands us $\alpha$

Solution:  $\alpha = \alpha_1, \dots, \alpha_N$

$$\Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

From the optimization:  
for  $n = 1, \dots, N$

$$\alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1) = 0$$

$$\alpha_n > 0 \Rightarrow y_n (\mathbf{w}^T \mathbf{x}_n + b) = 1 \Rightarrow \mathbf{x}_n \text{ is a } \boxed{\text{support vector}}$$

Intuition – remember we maximize  
w.r.t  $\alpha$

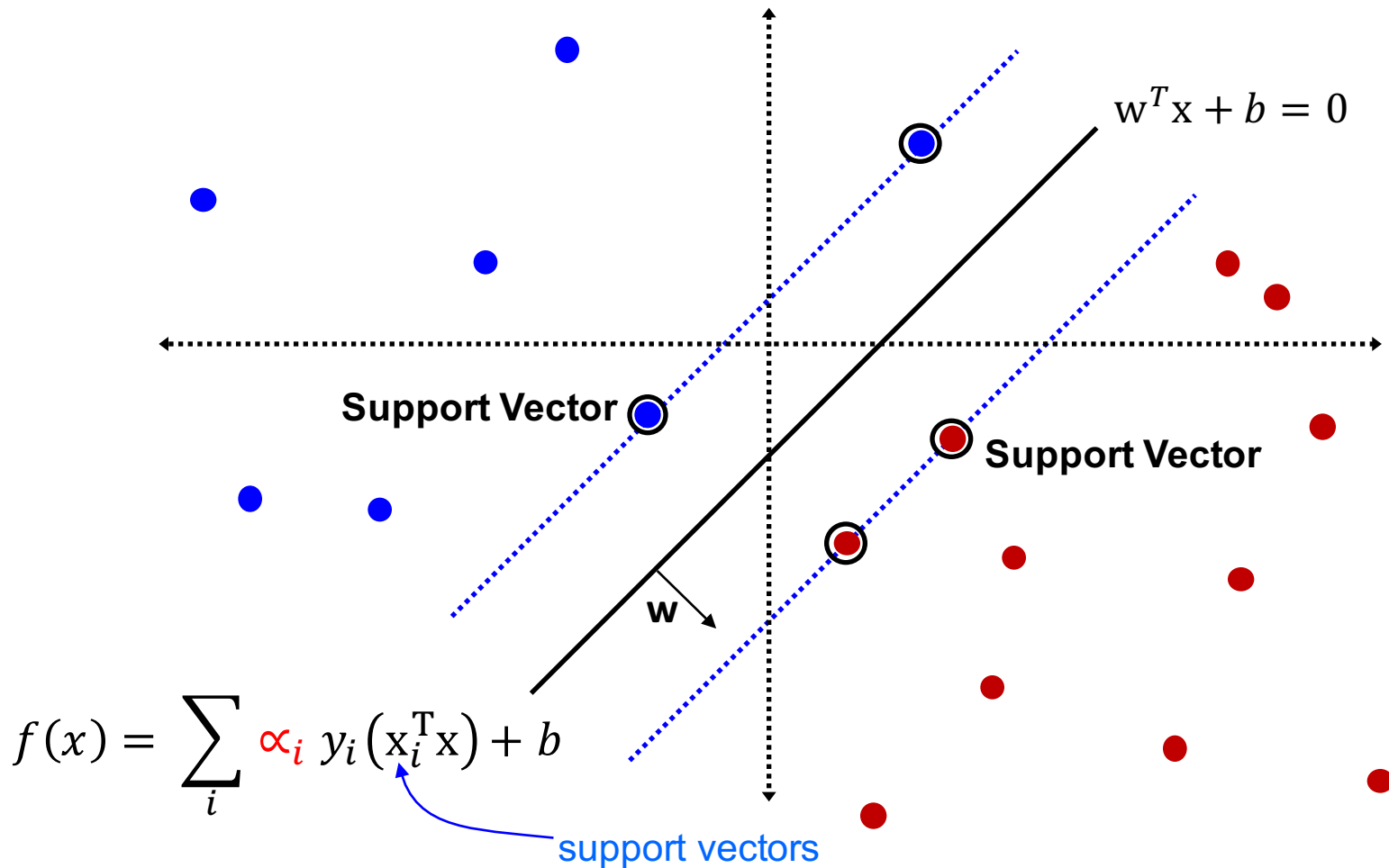
$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n (\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

How do we get  $b$ ?



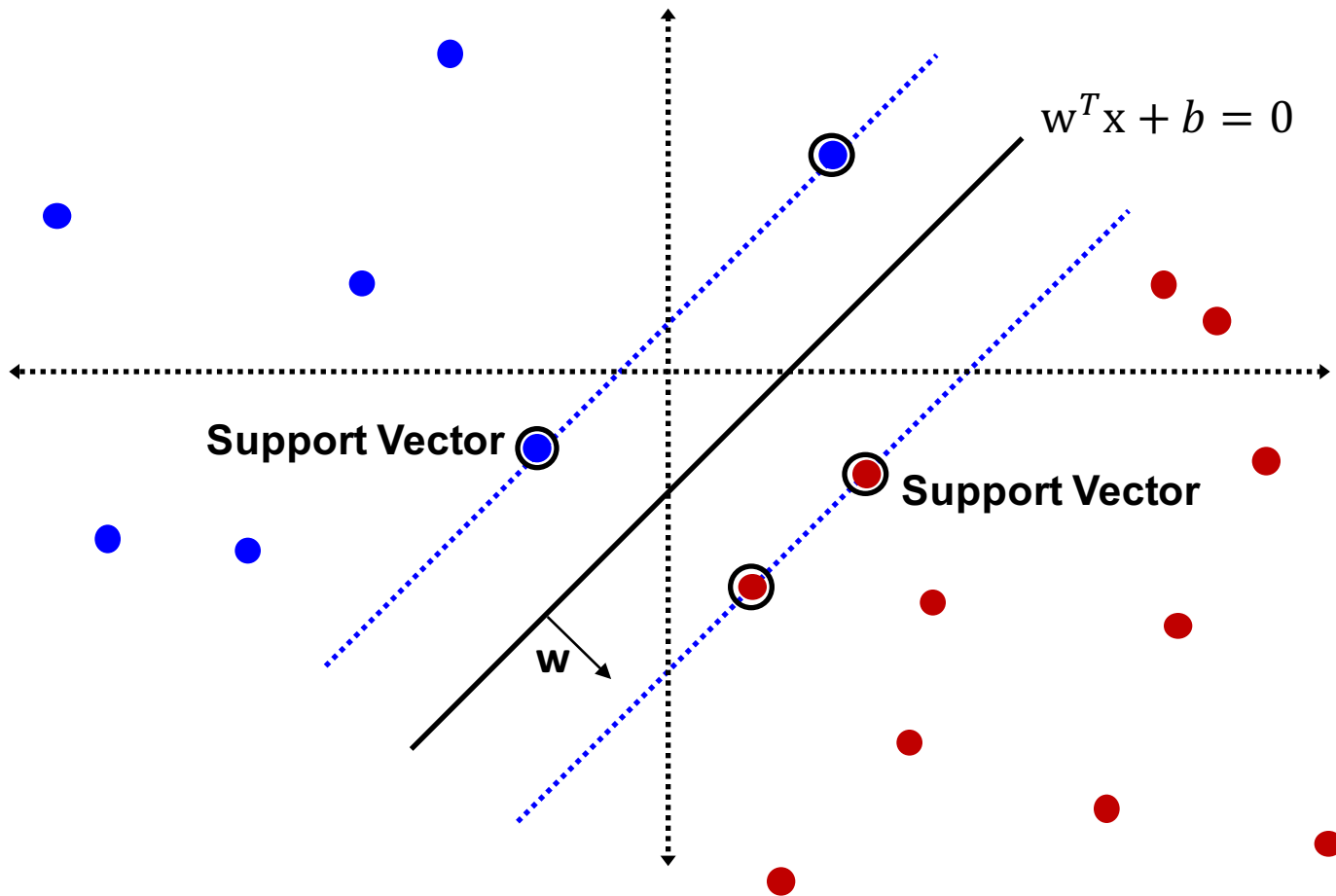
# Support Vector Machine

linearly separable data



# Leave One Out (LOO) error bound

How many LOO mistakes can we make in the worst case?

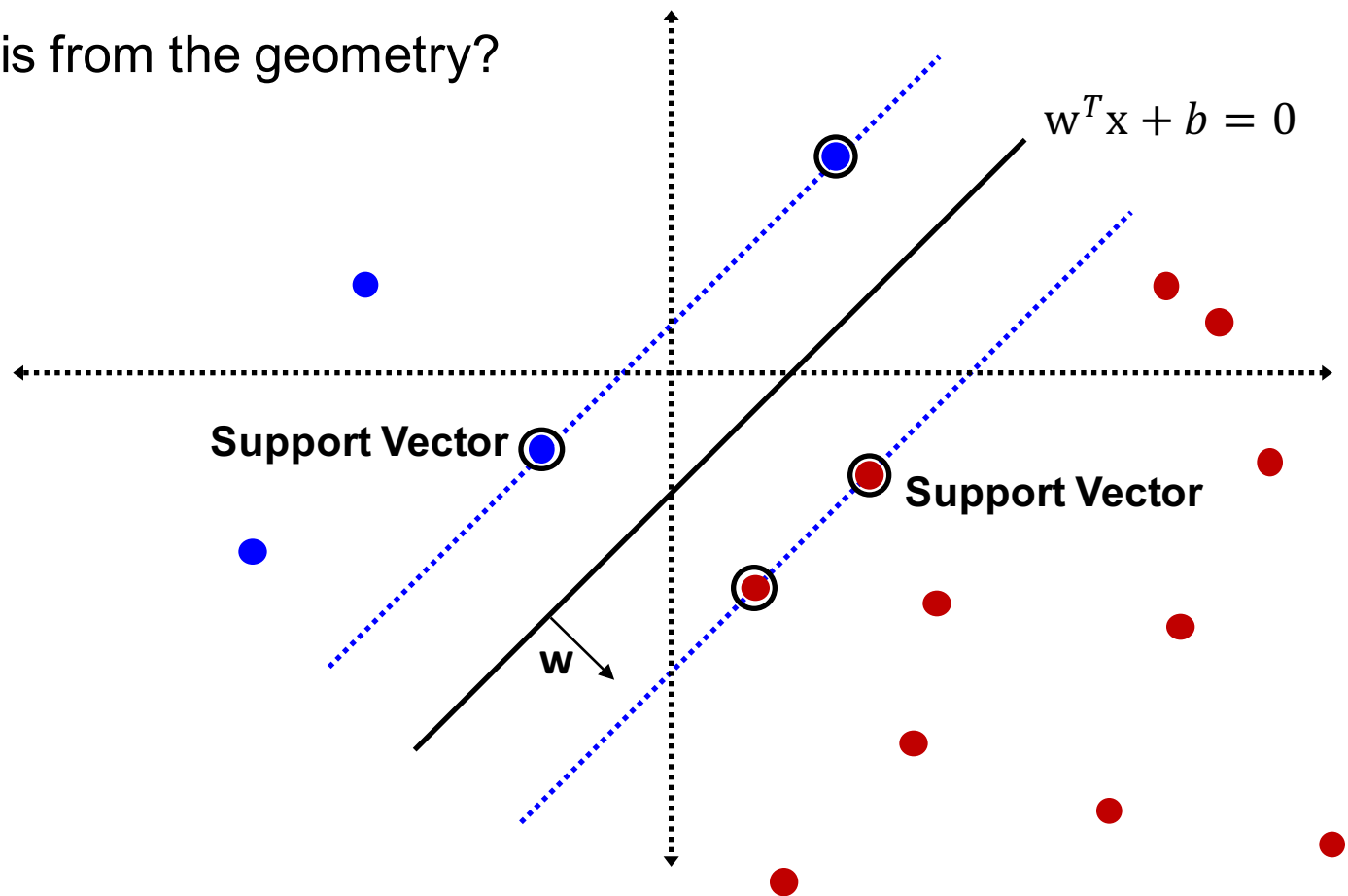


# Leave One Out (LOO) error bound

How many LOO mistakes can we make in the worst case?

#SV, i.e., the LOO error is bounded by #SV/N

Can you see this from the geometry?



# Leave One Out (LOO) error bound

How many LOO mistakes can we make in the worse case?

#SV, i.e., the LOO error is bounded by #SV/N

Can you see this from the primal QP?

$$\text{Minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{subject to } y_n (\mathbf{w}^T x_n + b) \geq 1 \\ n = 1, 2, \dots, N$$

$$\mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}$$

# Leave One Out (LOO) error bound

How many LOO mistakes can we make in the worse case?

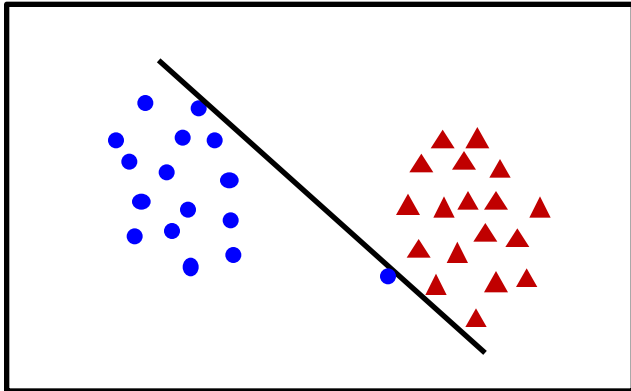
#SV, i.e., the LOO error is bounded by #SV/N

Can you see this from the dual QP?

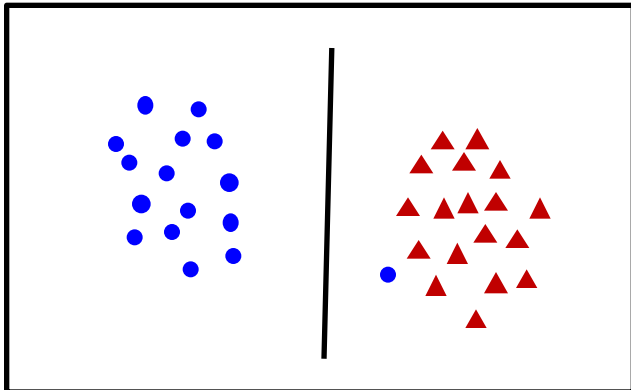
$$\min_{\alpha} \frac{1}{2} \alpha^T \begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \dots & y_1 y_N x_1^T x_N \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 & \dots & y_2 y_N x_2^T x_N \\ \dots & \dots & \dots & \dots \\ y_N y_1 x_N^T x_1 & y_N y_2 x_N^T x_2 & \dots & y_N y_N x_N^T x_N \end{bmatrix} \alpha + (-1^T) \alpha$$

$$\text{subject to } y^T \alpha = 0 \quad 0 \leq \alpha \leq \infty$$

# Linear separability again: What is the best $w$ ?



- linearly separated but very narrow margin

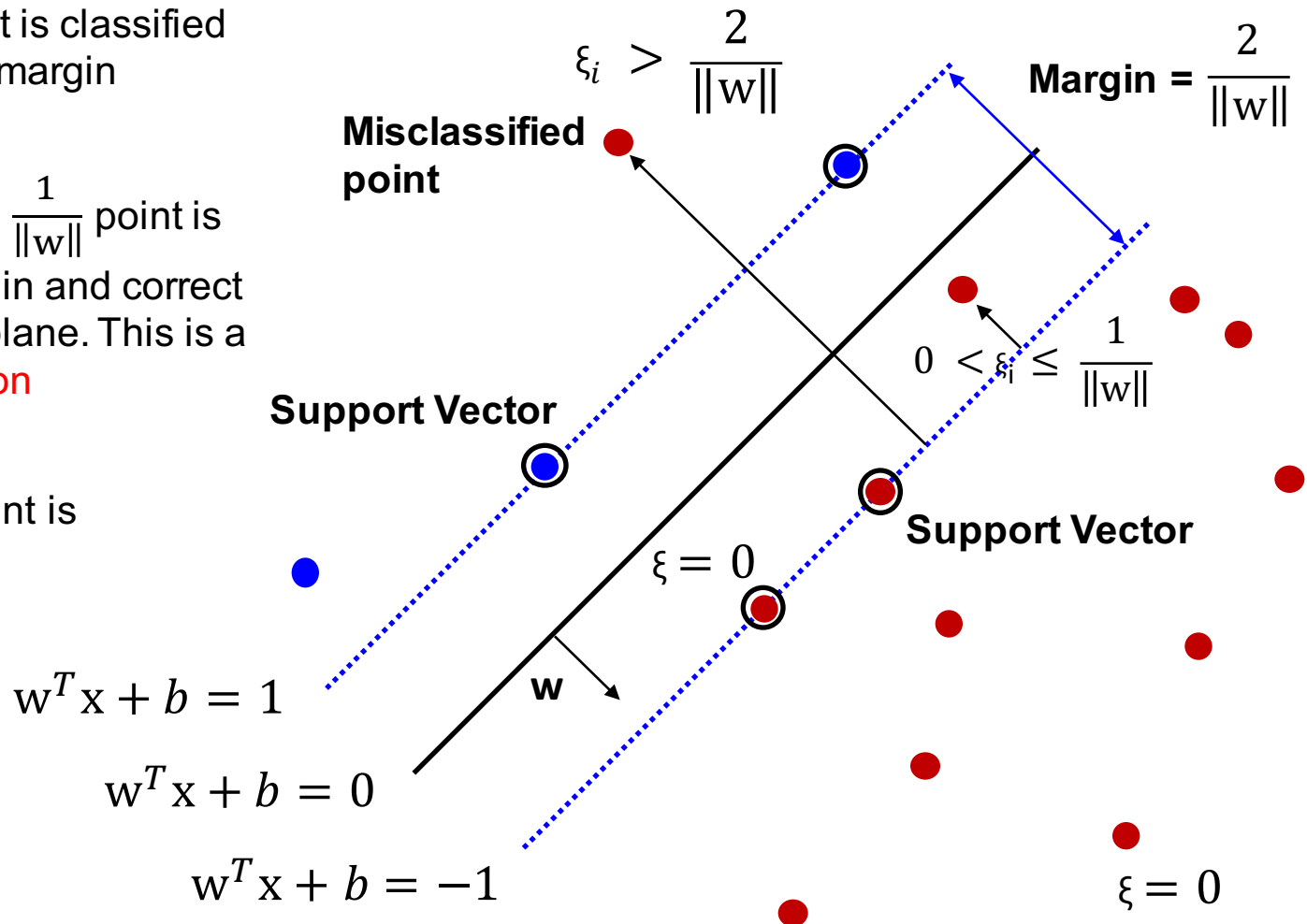


- possibly the large margin solution is better, even though one constraint is violated

- trade off between the margin and the number of mistakes on the training data

# Introduce “slack” variables $\xi_i \geq 0$

- For  $\xi = 0$  point is classified correctly with margin
- For  $0 < \xi \leq \frac{1}{\|w\|}$  point is between margin and correct side of hyperplane. This is a **margin violation**
- for  $\xi > \frac{1}{\|w\|}$  point is **misclassified**



# “Soft” margin solution

The PRIMAL optimization problem becomes the following QP problem

$$\min_{w \in \mathbb{R}^d, \xi_i > 0} \|w\|^2 + C \sum_i^N \xi_i$$

subject to

$$y_i (w^T x_i + b) \geq 1 - \xi_i \text{ for } i = 1 \dots N$$

- Every constraint can be satisfied if  $\xi_i$  is sufficiently large
- $C$  is a **regularization** parameter:
  - small  $C$  allows constraints to be easily ignored  $\rightarrow$  large margin
  - large  $C$  makes constraints hard to ignore  $\rightarrow$  narrow margin
  - $C = \infty$  enforces all constraints: hard margin
- There is still a unique minimum
- Pain: we have now a parameter  $C$
- Pain2: adding apples and oranges



# The dual formulation

$$\min_{\alpha} \frac{1}{2} \alpha^T \underbrace{\begin{bmatrix} y_1 y_1 x_1^T x_1 & y_1 y_2 x_1^T x_2 & \dots & y_1 y_N x_1^T x_N \\ y_2 y_1 x_2^T x_1 & y_2 y_2 x_2^T x_2 & \dots & y_2 y_N x_2^T x_N \\ \dots & \dots & \dots & \dots \\ y_N y_1 x_N^T x_1 & y_N y_2 x_N^T x_2 & \dots & y_N y_N x_N^T x_N \end{bmatrix}}_{\text{quadratic coefficients}} \alpha + \underbrace{(-1^T)}_{\text{linear}} \alpha$$

subject to  $\underbrace{y^T \alpha}_{\text{linear constraint}} = 0$

$$\underbrace{0}_{\text{lower bounds}} \leq \alpha \leq \underbrace{C}_{\text{upper bounds}}$$

# Another point of view for SVM

SVM was formulated as a **constrained** optimization problem over  $w$  and positive  $\xi$

$$\min_{w \in \mathbb{R}^d, \xi_i > 0} \|w\|^2 + C \sum_i^N \xi_i \quad \text{subject to} \quad y_i(w^T x_i + b) \geq 1 - \xi_i \text{ for } i = 1 \dots N$$

The constraint  $y_i(w^T x_i + b) \geq 1 - \xi_i$  can be written as

$$y_i f(x_i) \geq 1 - \xi_i$$

since  $\xi_i \geq 0$  and since we want to “pay” as little as possible

$$\xi_i = \max(0, 1 - y_i f(x_i))$$

Substituting, we get the **unconstrained** optimization problem over  $w$  alone

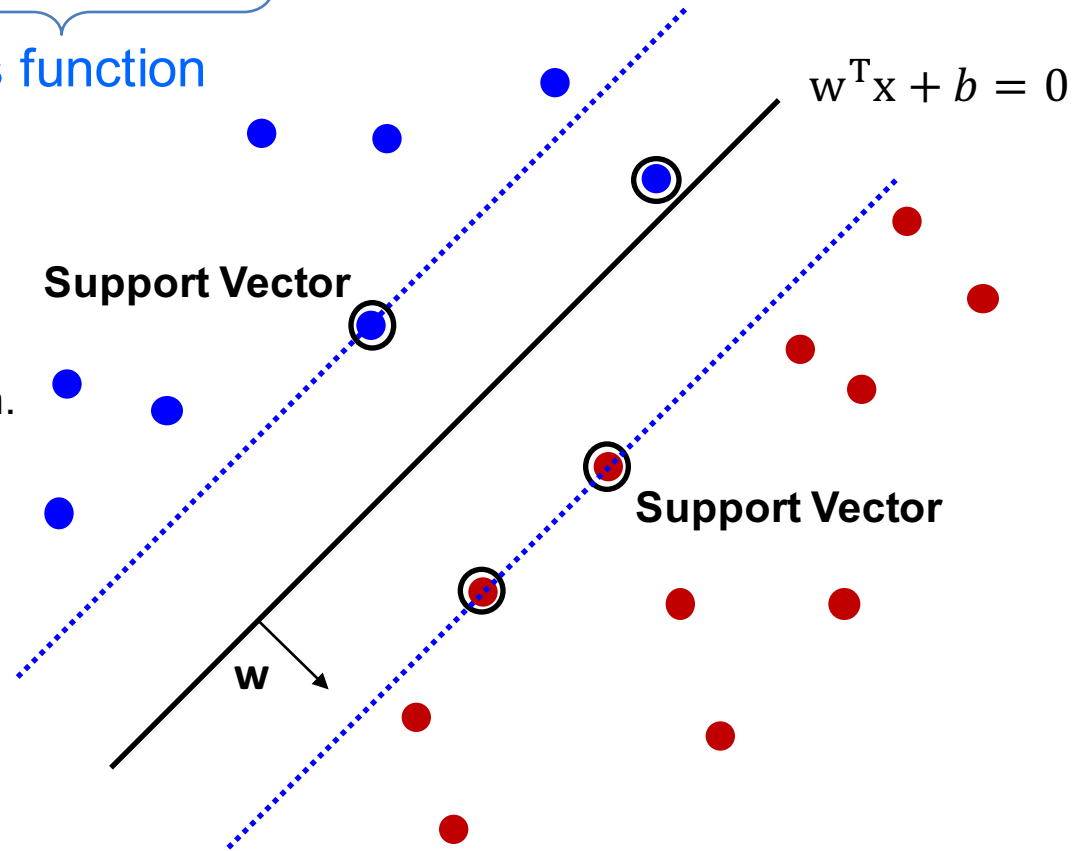
$$\min_{w \in \mathbb{R}^d} \underbrace{\|w\|^2}_{\text{regularization}} + C \sum_i^N \underbrace{\max(0, 1 - y_i f(x_i))}_{\text{loss function}}$$

# Loss function

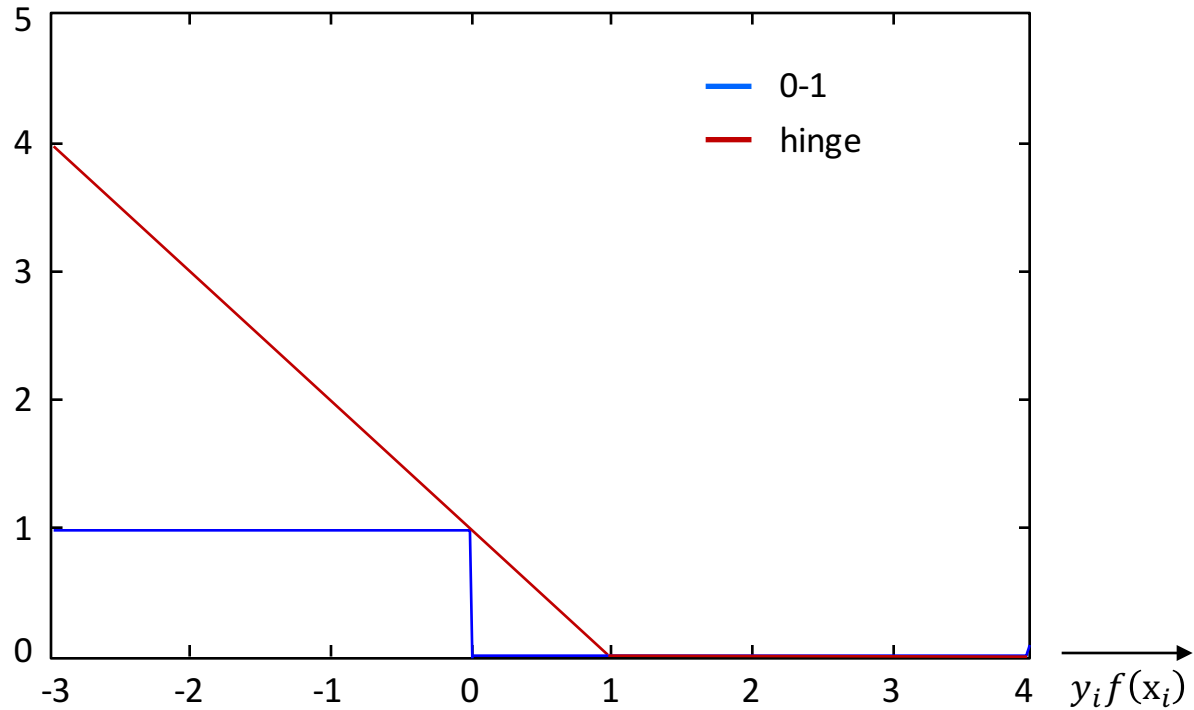
$$\min_{w \in \mathbb{R}^d} \|w\|^2 + C \sum_i^N \underbrace{\max(0, 1 - y_i f(x_i))}_{\text{loss function}}$$

Points are in two categories:

1.  $y_i f(x_i) \geq 1$   
Point is on or outside margin.  
No contribution to loss
2.  $y_i f(x_i) < 1$   
Point violates margin  
constraint. Contributes to  
loss.



# Hinge loss



- SVM uses the hinge loss  $\max(0, 1 - y_i f(x_i))$
- an approximation to the 0-1 loss

# What about multi-class SVMs?

- Unfortunately, there is no “definitive” multi-class SVM formulation
- In practice, we have to obtain a multi-class SVM by combining multiple two-class SVMs
- One vs. others
  - Training: learn an SVM for each class vs. the others
  - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
- One vs. one
  - Training: learn an SVM for each pair of classes
  - Testing: each learned SVM “votes” for a class to assign to the test example

# SVM: Practical Advice

- Avoid MATLAB's implementation
- Use <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
- Or <http://www.csie.ntu.edu.tw/~cjlin/liblinear/>  
when you know you want linear SVM
- Play with C (?)
- Avoid unbalanced classes (?)

# SVM – sketch derivation

- Since  $w^T x + b = 0$  and  $c(w^T x + b) = 0$  define the same plane, we have the freedom to choose the normalization of  $w$
- Choose normalization such that  $w^T x_+ + b = +1$  and  $w^T x_- + b = -1$  for the positive and negative support vectors respectively
- Then the **margin** is given by

$$\frac{w}{\|w\|} \cdot (x_+ - x_-) = \frac{w^T (x_+ - x_-)}{\|w\|} = \frac{2}{\|w\|}$$

# SVM – Optimization

- Learning the SVM can be formulated as an optimization:

$$\max_w \frac{2}{\|w\|} \quad \text{subject to } w^T x_i + b \begin{cases} \geq 1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \quad \text{for } i = 1 \dots N$$

- Or equivalently

$$\min_w \|w\|^2 \quad \text{subject to } y_i (w^T x_i + b) \geq 1 \quad \text{for } i = 1 \dots N$$

- This is a quadratic optimization problem subject to linear constraints and there is a unique minimum



# SVM – sketch derivation

- Since  $w^T x + b = 0$  and  $c(w^T x + b) = 0$  define the same plane, we have the freedom to choose the normalization of  $w$
- Choose normalization such that  $w^T x_+ + b = +1$  and  $w^T x_- + b = -1$  for the positive and negative support vectors respectively
- Then the **margin** is given by

$$\frac{w}{\|w\|} \cdot (x_+ - x_-) = \frac{w^T (x_+ - x_-)}{\|w\|} = \frac{2}{\|w\|}$$

# Finding $w$ with large margin

Let  $x_n$  be the nearest data point to the plane  $w^T x = 0$ . How far is it?

2 preliminary technicalities:

1. **Normalize  $w$ :**

$$|w^T x_n| = 1$$

2. **Pull out  $\omega_0$ :**

$$w = (\omega_1, \dots, \omega_d) \text{ apart from } b$$

The plane is now  $w^T x + b = 0$  (no  $x_0$ )