

## Linear classifier (E.g., Perceptron)

$$h(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

Which outputs  $+1$  or  $-1$ .

Say:

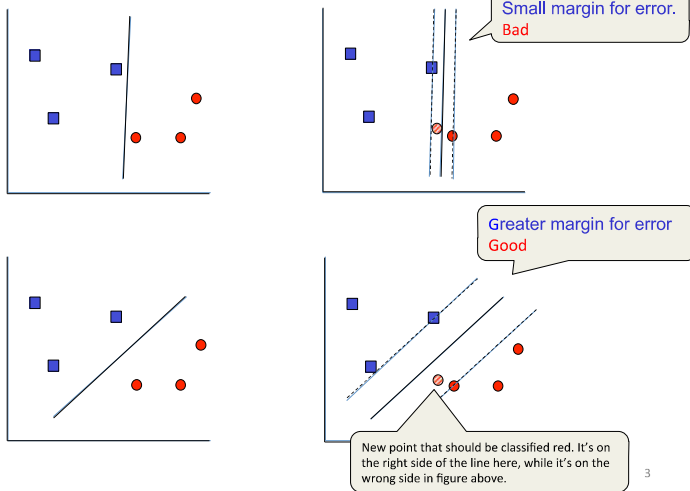
$+1$  corresponds to blue, and  
 $-1$  to red, or vice versa.

Many lines do the job.  
 Which one to choose?

## Support Vector Machines

Many of these slides are derived from Seyong Kim and Alex Thomo. Thanks!

### Margin for different separators



### Scale Invariance

$$h(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

- We rescale  $\mathbf{w}$  and  $b$  (without changing the line) such that:

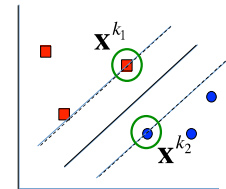
$$\mathbf{w} \cdot \mathbf{x}^{k_1} + b = 1$$

for the closest point(s) to the line on the  $+1$  side, and

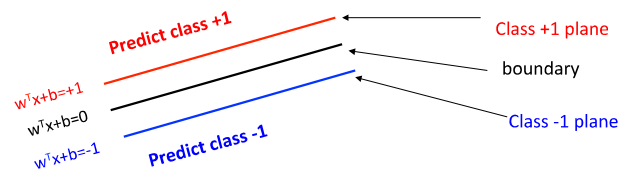
$$\mathbf{w} \cdot \mathbf{x}^{k_2} + b = -1$$

for the closest point(s) to the line on the  $-1$  side.

Closest points are called "support vectors".

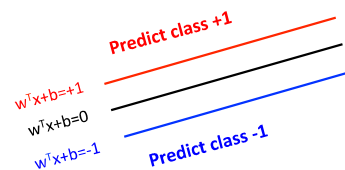


### Specifying a max margin classifier



Classify as $+1$	if	$\mathbf{w}^T \mathbf{x} + b \geq 1$
Classify as $-1$	if	$\mathbf{w}^T \mathbf{x} + b \leq -1$
Undefined	if	$-1 < \mathbf{w}^T \mathbf{x} + b < 1$

### Specifying a max margin classifier

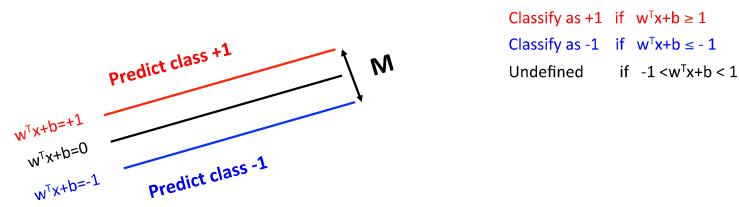


Classify as $+1$	if	$\mathbf{w}^T \mathbf{x} + b \geq 1$
Classify as $-1$	if	$\mathbf{w}^T \mathbf{x} + b \leq -1$
Undefined	if	$-1 < \mathbf{w}^T \mathbf{x} + b < 1$

Is the linear separation assumption realistic?

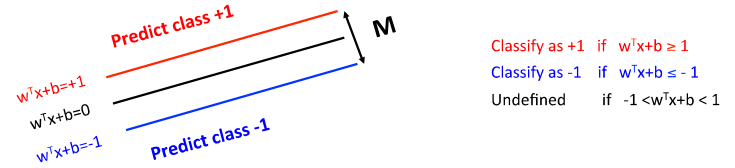
We will deal with this shortly, but let's assume it for now

## Maximizing the margin



- Lets define the width of the margin by  $M$
- How can we encode our goal of maximizing  $M$  in terms of our parameters ( $w$  and  $b$ )?
- Lets start with a few observations

## Maximizing the margin



- Observation 1: the vector  $w$  is orthogonal to the +1 plane
- Why?

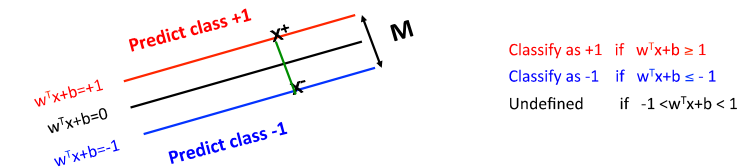
Let  $u$  and  $v$  be two points on the +1 plane, then for the vector defined by  $u$  and  $v$  we have  $w^T(u-v) = 0$

Corollary: the vector  $w$  is orthogonal to the -1 plane

7

8

## Maximizing the margin



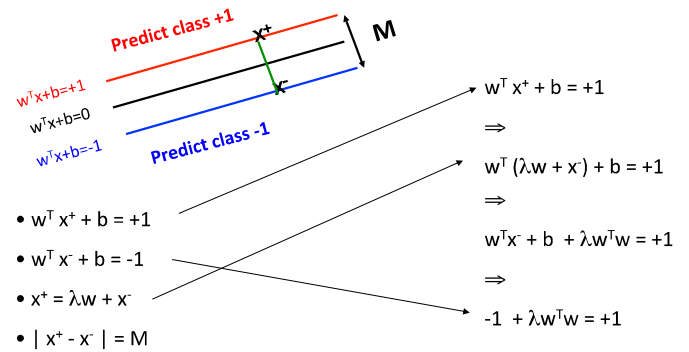
- Observation 1: the vector  $w$  is orthogonal to the +1 and -1 planes
- Observation 2: if  $x^+$  is a point on the +1 plane and  $x^-$  is the closest point to  $x^+$  on the -1 plane then

$$x^+ = \lambda w + x^-$$

Since  $w$  is orthogonal to both planes we need to 'travel' some distance along  $w$  to get from  $x^+$  to  $x^-$

9

## Putting it together

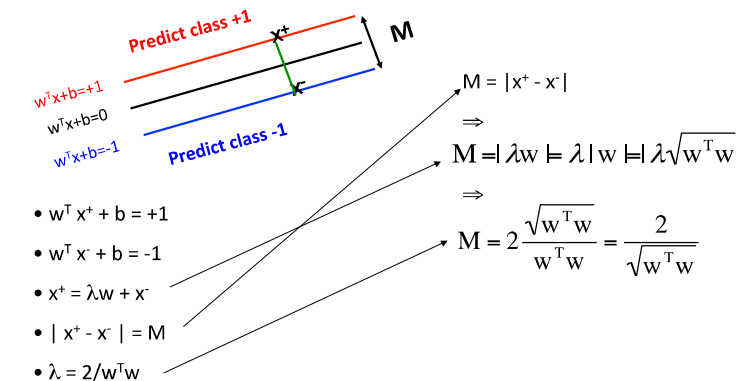


We can now define  $M$  in terms of  $w$  and  $b$

$$\lambda = 2/w^T w$$

10

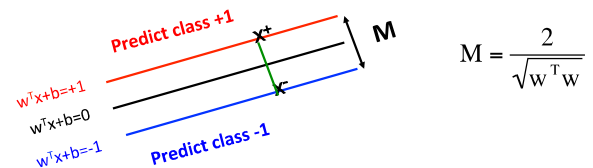
## Putting it together



We can now define  $M$  in terms of  $w$  and  $b$

11

## Finding the optimal parameters



We can now search for the optimal parameters by finding a solution that:

1. Correctly classifies all points
2. Maximizes the margin (or equivalently minimizes  $w^T w$ )

Several optimization methods can be used: Gradient descent, simulated annealing, EM etc.

12

## Aside: Quadratic programming (QP)

Quadratic programming solves optimization problems of the following form:

$$\min_u \frac{u^T R u}{2} + d^T u + c$$

subject to  $n$  inequality constraints:

$$a_{11}u_1 + a_{12}u_2 + \dots \leq b_1$$

$$M \quad M \quad M$$

$$a_{n1}u_1 + a_{n2}u_2 + \dots \leq b_n$$

and  $k$  equality constraints:

$$a_{n+1,1}u_1 + a_{n+1,2}u_2 + \dots = b_{n+1}$$

$$M \quad M \quad M$$

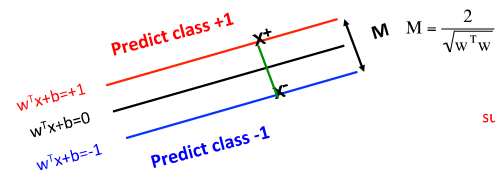
$$a_{n+k,1}u_1 + a_{n+k,2}u_2 + \dots = b_{n+k}$$

Quadratic term

When a problem can be specified as a QP problem we can use solvers that are better than gradient descent or simulated annealing

13

## SVM as a QP problem



$$\min_u \frac{u^T R u}{2} + d^T u + c$$

subject to  $n$  inequality constraints:

$$a_{11}u_1 + a_{12}u_2 + \dots \leq b_1$$

$$M \quad M \quad M$$

$$a_{n1}u_1 + a_{n2}u_2 + \dots \leq b_n$$

and  $k$  equality constraints:

$$a_{n+1,1}u_1 + a_{n+1,2}u_2 + \dots = b_{n+1}$$

$$M \quad M \quad M$$

$$a_{n+k,1}u_1 + a_{n+k,2}u_2 + \dots = b_{n+k}$$

$$\text{Min } (w^T w)/2$$

subject to the following inequality constraints:

For all  $x$  in class +1

$$w^T x + b \geq 1$$

For all  $x$  in class -1

$$w^T x + b \leq -1$$

A total of  $n$  constraints if we have  $n$  input samples

14

## SVM as a QP problem: a simplification

$$\text{Min } (w^T w)/2$$

subject to the following inequality constraints:

For all  $x$  in class +1

$$w^T x + b \geq 1$$

For all  $x$  in class -1

$$w^T x + b \leq -1$$

$$\text{Min } (w^T w)/2$$

subject to the following inequality constraints:

For all  $x$  in class +1

$$y(w^T x + b) \geq 1$$

For all  $x$  in class -1

$$y(w^T x + b) \geq 1$$

The same constraint!!  
So much easier to handle!

15

## Example

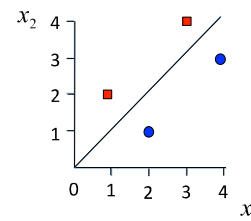
Training tuples:

$([1, 2], +1)$

$([2, 1], -1)$

$([3, 4], +1)$

$([4, 3], -1)$



$$\min_{w_1, w_2} \frac{1}{2} (w_1^2 + w_2^2)$$

subject to

$$(+1)(w_1 + 2w_2) \geq 1$$

$$(-1)(2w_1 + w_2) \geq 1$$

$$(+1)(3w_1 + 4w_2) \geq 1$$

$$(-1)(4w_1 + 3w_2) \geq 1$$

For this example, solution is easy to see:

$b = 0$  and  $w = [-1, +1]$ , i.e. the line is:

$$-x_1 + x_2 = 0$$

All conditions are satisfied.

$$M = \frac{2}{\|w\|} = \frac{2}{\sqrt{1+1}} = \frac{2}{\sqrt{2}} = \sqrt{2}$$

16

## Non linearly separable case

• So far we assumed that a linear plane can perfectly separate the points

• But this is not usually the case

- noise, outliers

How can we convert this to a QP problem?

- Minimize training errors?

$$\min w^T w$$

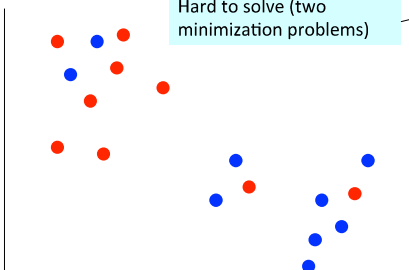
$$\min \text{\#errors}$$

- Penalize training errors:

$$\min w^T w + C * (\text{\#errors})$$

Hard to solve (two minimization problems)

Hard to encode in a QP problem



17

## Non linearly separable case

• Instead of minimizing the number of misclassified points we can minimize the distance between these points and their correct plane

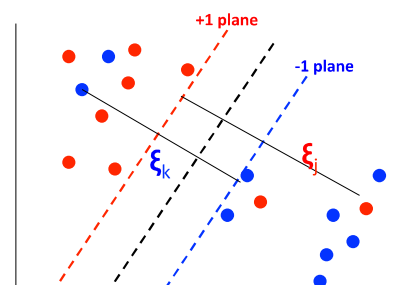
The new optimization problem is:

$$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \xi_i$$

subject to the following inequality constraints:

For all  $x_i$  in class +1 or class -1

$$y^*(w^T x + b) \geq 1 - \xi_i$$



Wait. Are we missing something?

18

## Final optimization for non linearly separable case

The new optimization problem is:

$$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \xi_i$$

subject to the following inequality constraints:

For all  $x_i$  in class +1

or class -1

$$y^*(w^T x + b) \geq 1 - \xi_i$$

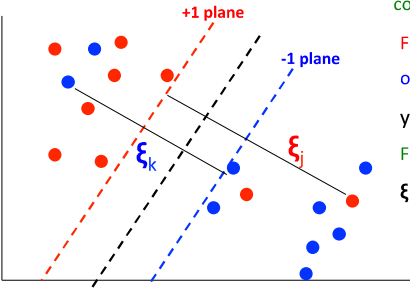
For all  $i$

$$\xi_i \geq 0$$

n constraints

Another n constraints

Slack variables ( $\epsilon$ )



19

## Where we are

Two optimization problems: For the separable and non separable cases

$$\min_w \frac{w^T w}{2}$$

For all  $x$  in class + 1

$$w^T x + b \geq 1$$

For all  $x$  in class - 1

$$w^T x + b \leq -1$$

$$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \xi_i$$

For all  $x_i$  in class + 1

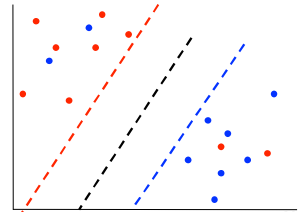
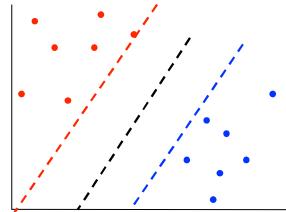
$$w^T x + b \geq 1 - \xi_i$$

For all  $x_i$  in class - 1

$$w^T x + b \leq -1 + \xi_i$$

For all  $i$

$$\xi_i \geq 0$$



20

## Where we are

Two optimization problems: For the separable and non separable cases

$$\min_w \frac{w^T w}{2}$$

For all  $x$  in class + 1

$$w^T x + b \geq 1$$

For all  $x$  in class - 1

$$w^T x + b \leq -1$$

$$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \xi_i$$

For all  $x_i$  in class + 1

$$w^T x + b \geq 1 - \xi_i$$

For all  $x_i$  in class - 1

$$w^T x + b \leq -1 + \xi_i$$

For all  $i$

$$\xi_i \geq 0$$

- Instead of solving these QPs directly we will solve a dual formulation of the SVM optimization problem

- The main reason for switching to this type of representation is that it would allow us to use a neat trick that will make our lives easier (and the run time faster)

21