

# Support Vector Machines

**Andrew W. Moore**

**Professor**

**School of Computer Science**

**Carnegie Mellon University**

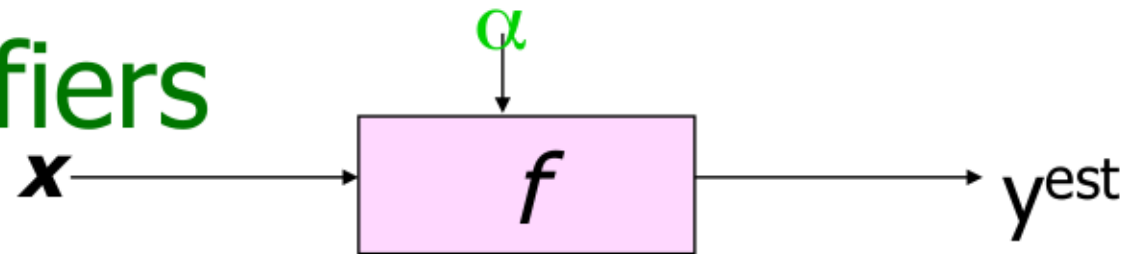
[www.cs.cmu.edu/~awm](http://www.cs.cmu.edu/~awm)

[awm@cs.cmu.edu](mailto:awm@cs.cmu.edu)

412-268-7599

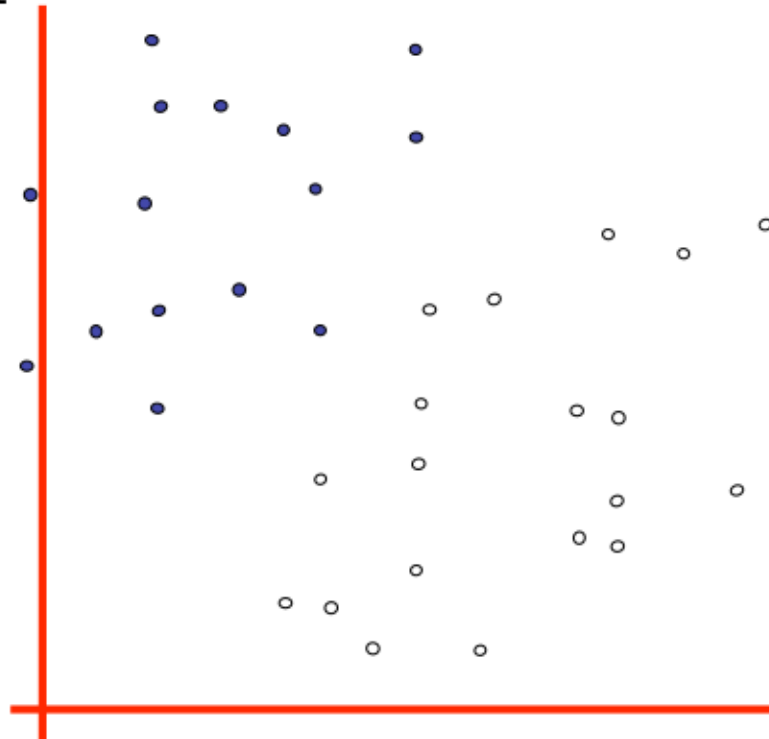
Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials>. Comments and corrections gratefully received.

# Linear Classifiers



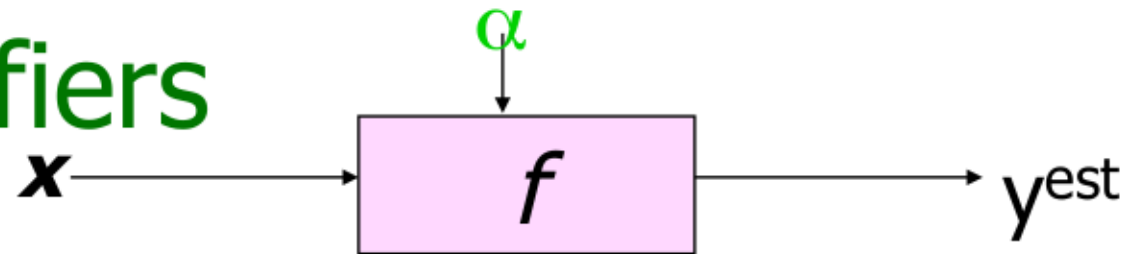
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

- denotes +1
- denotes -1

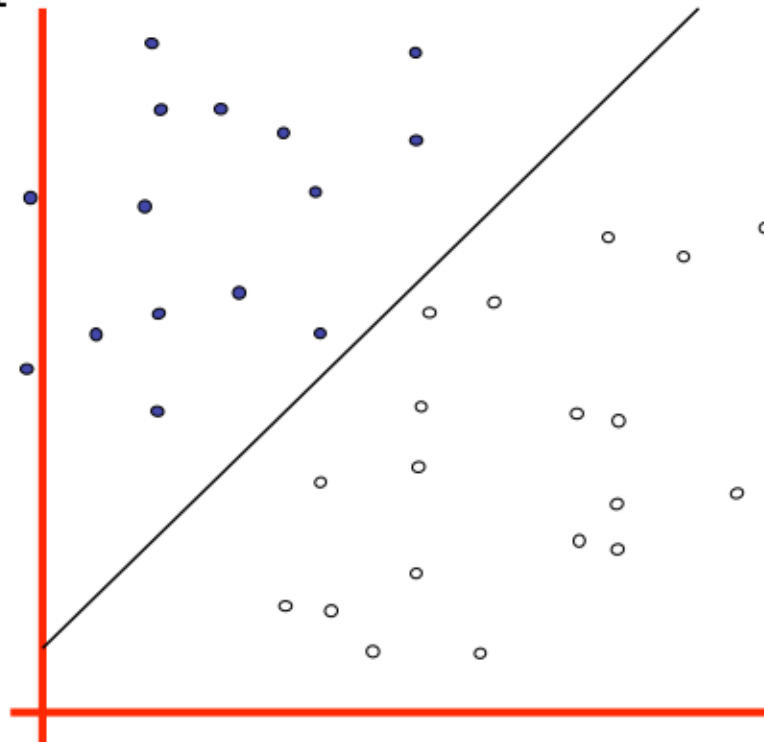


How would you classify this data?

# Linear Classifiers



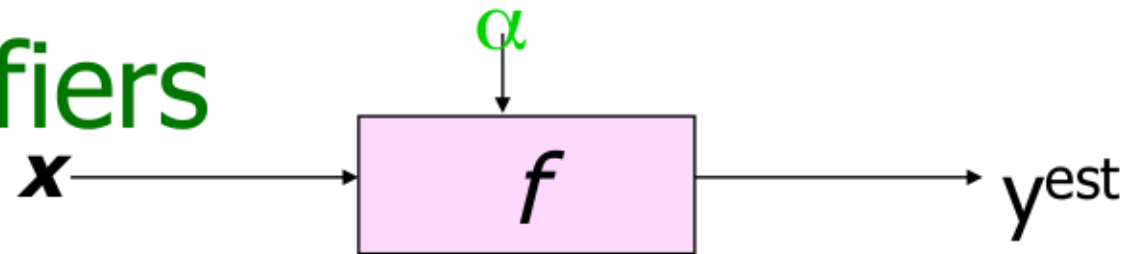
- denotes +1
- denotes -1



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

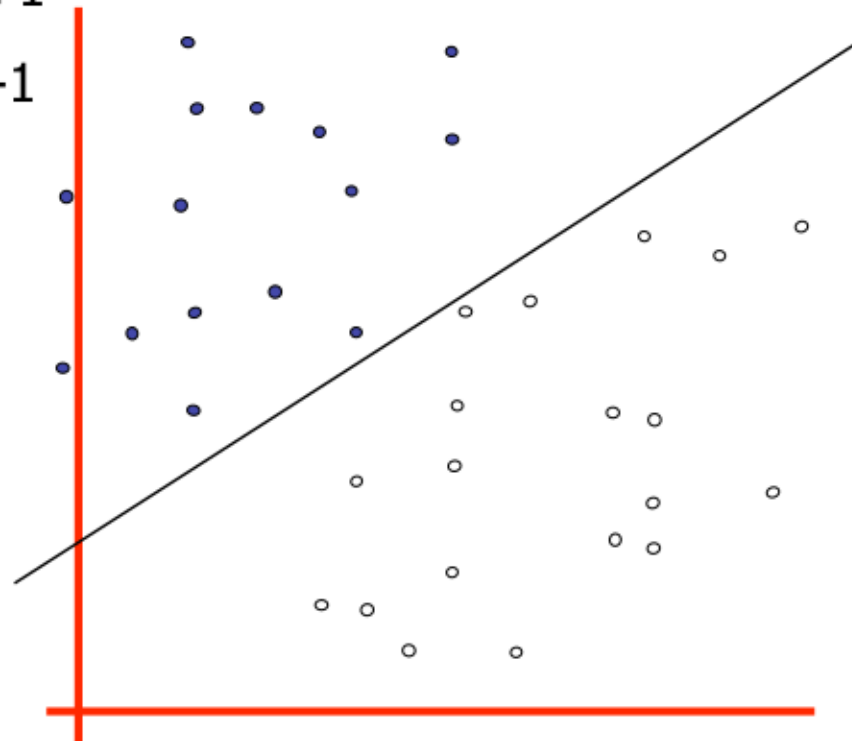
How would you classify this data?

# Linear Classifiers



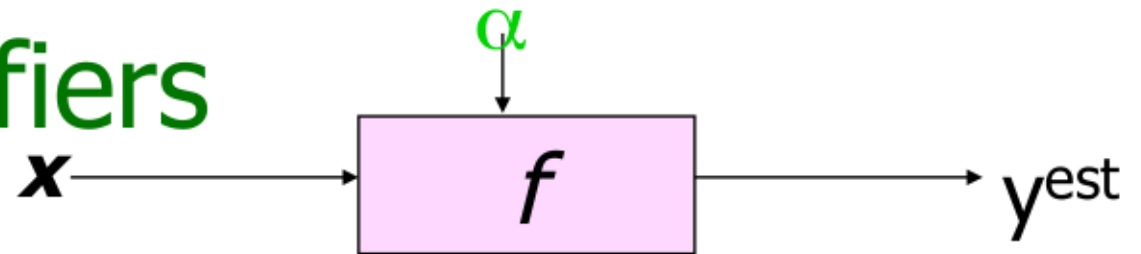
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

- denotes +1
- denotes -1

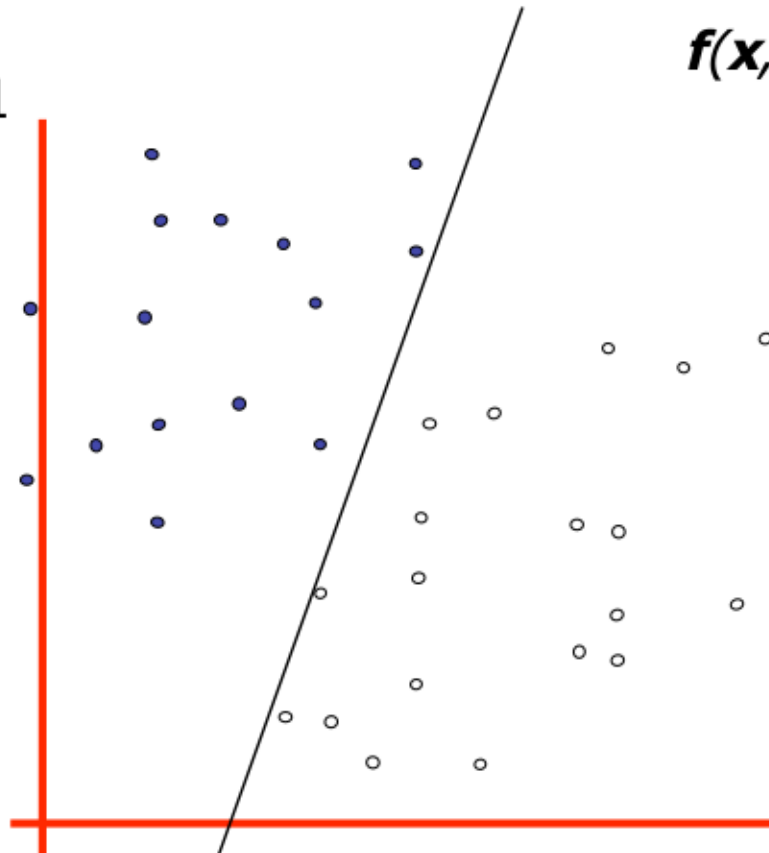


How would you classify this data?

# Linear Classifiers



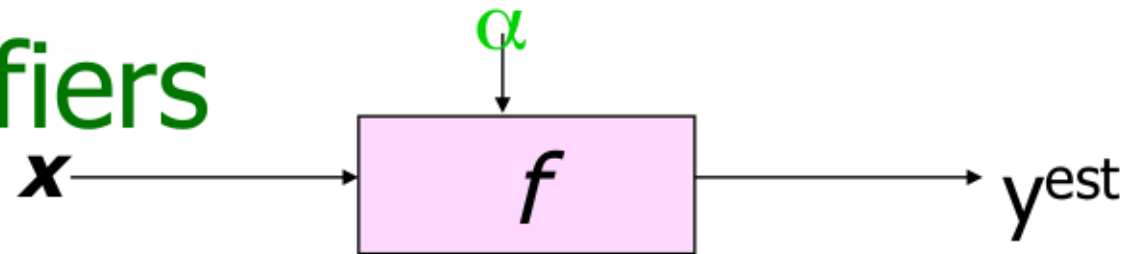
- denotes +1
- denotes -1



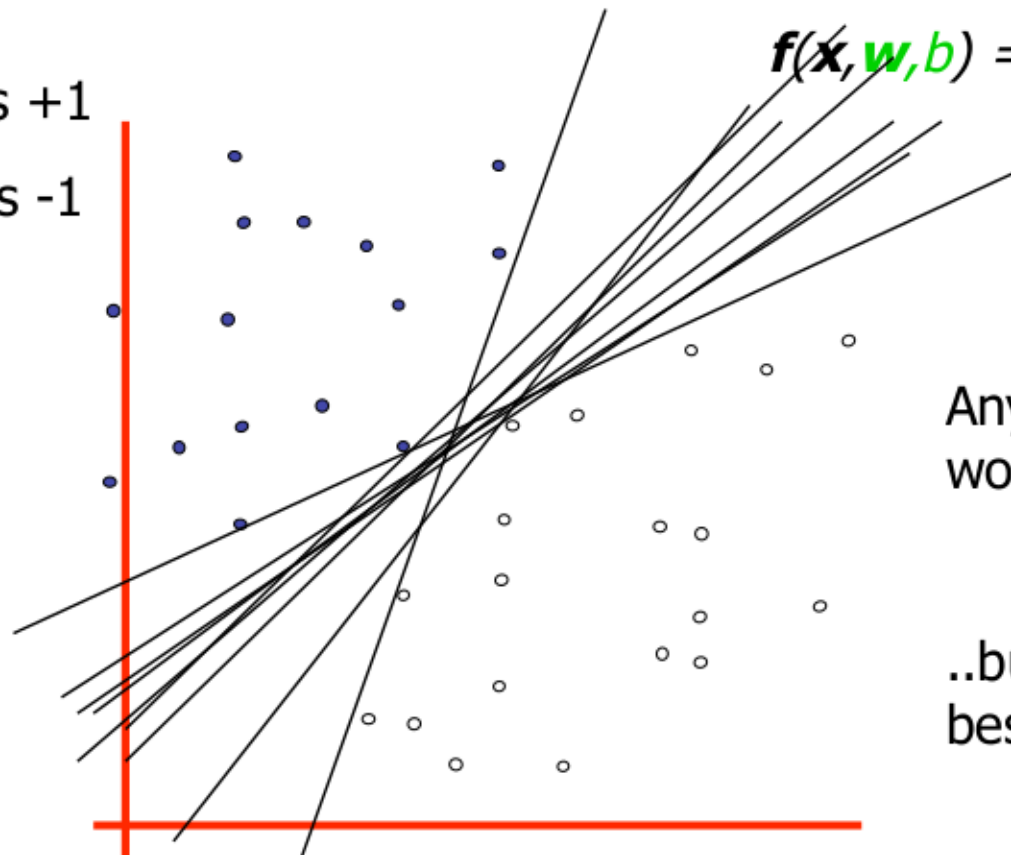
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

How would you classify this data?

# Linear Classifiers



- denotes +1
- denotes -1

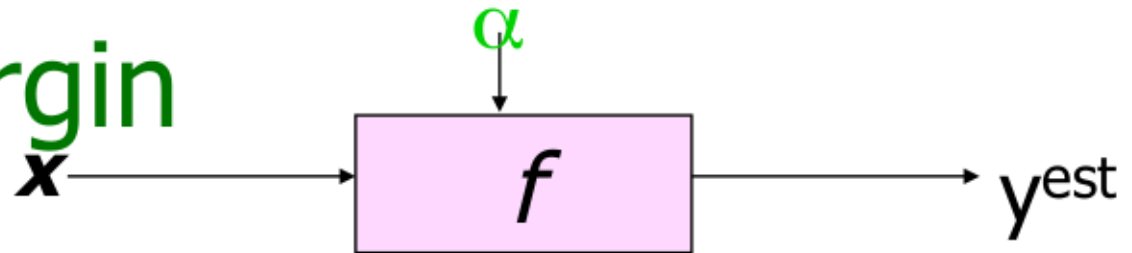


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

Any of these  
would be fine..

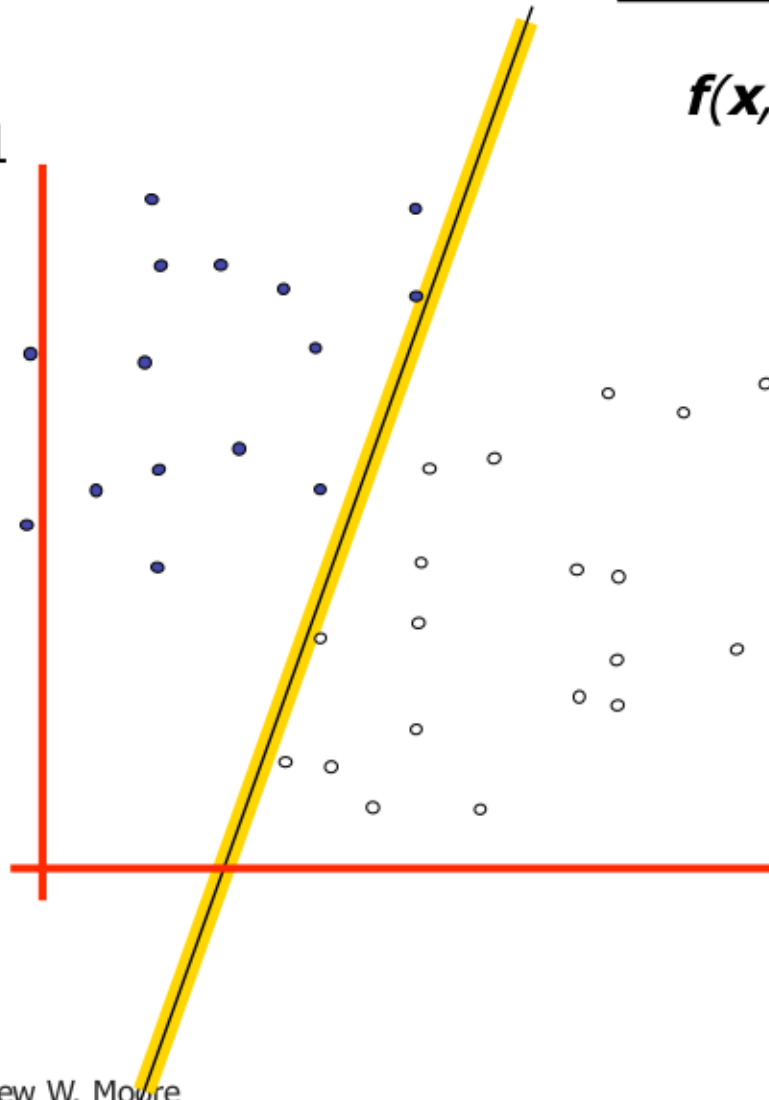
..but which is  
best?

# Classifier Margin



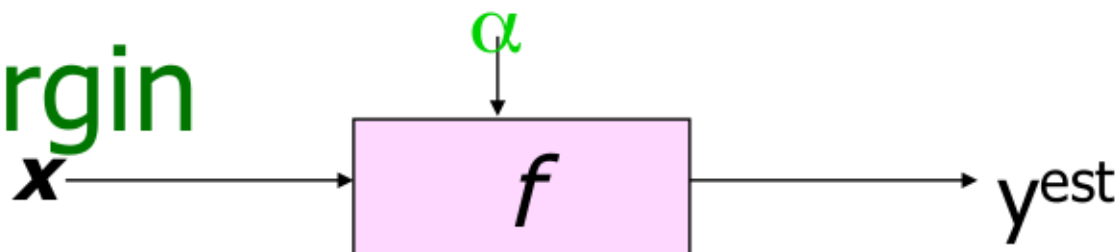
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

- denotes +1
- denotes -1

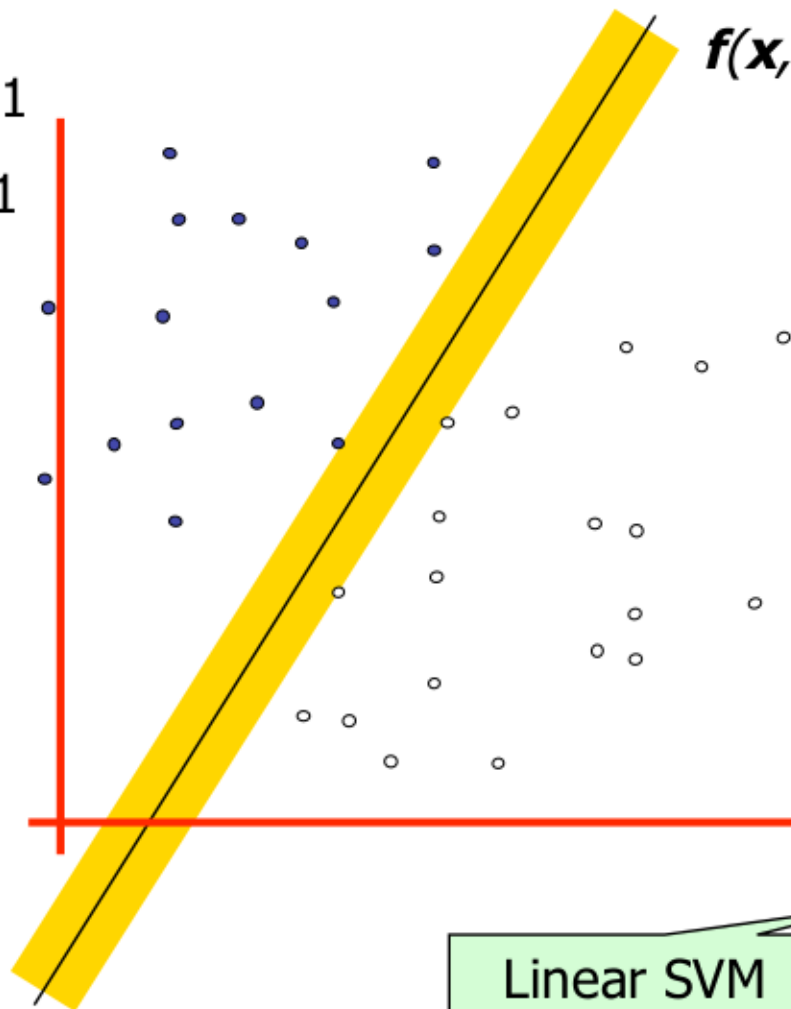


Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin



- denotes +1
- denotes -1



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

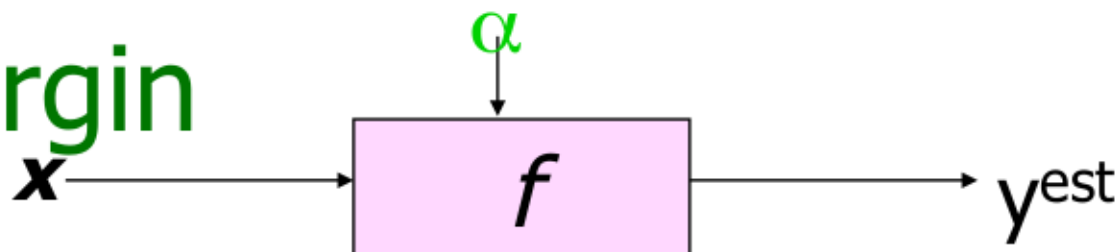
The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

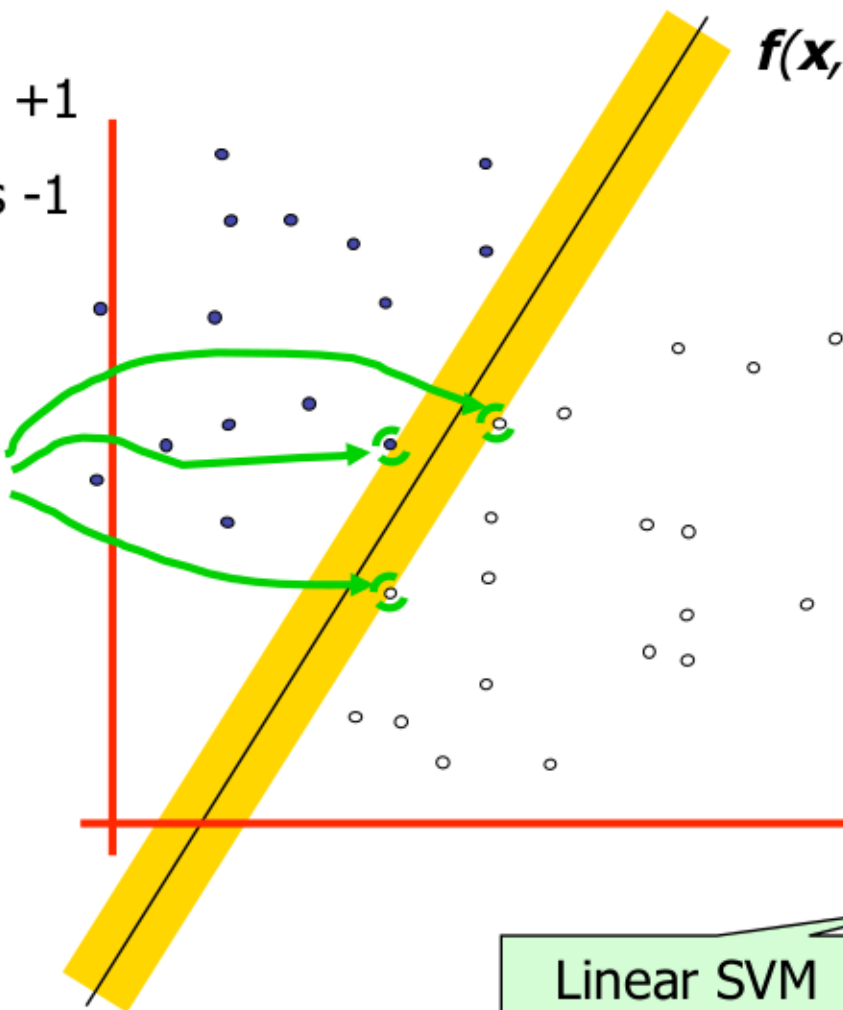


# Maximum Margin



- denotes +1
- denotes -1

**Support Vectors**  
are those  
datapoints that  
the margin  
pushes up  
against



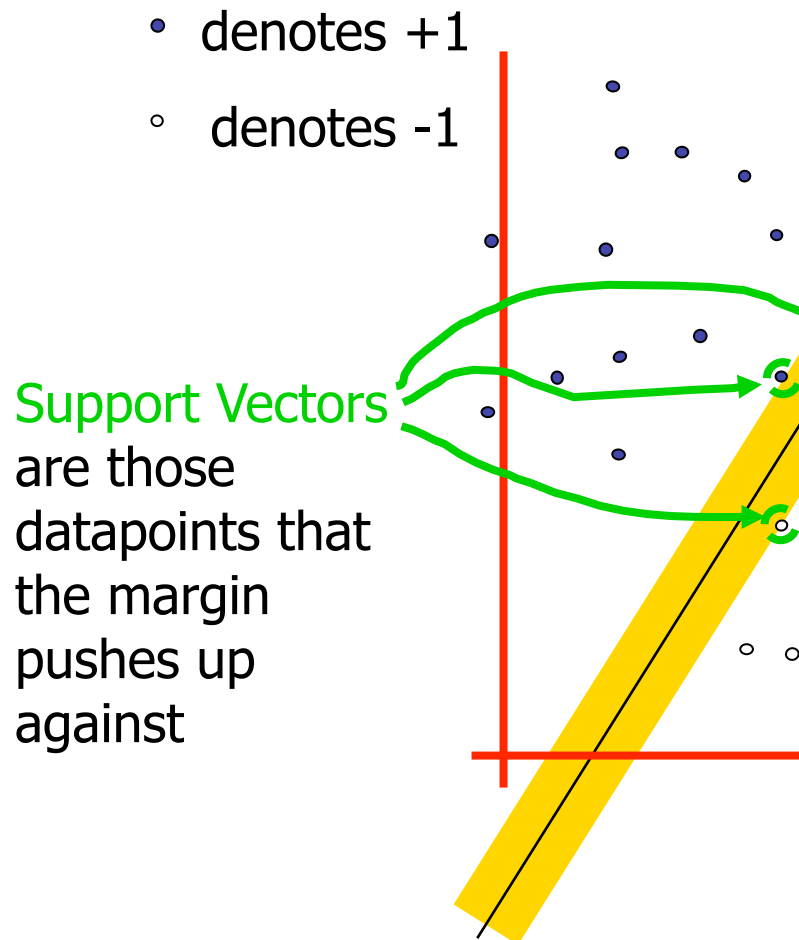
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin.

This is the simplest kind of SVM (Called an LSVM)

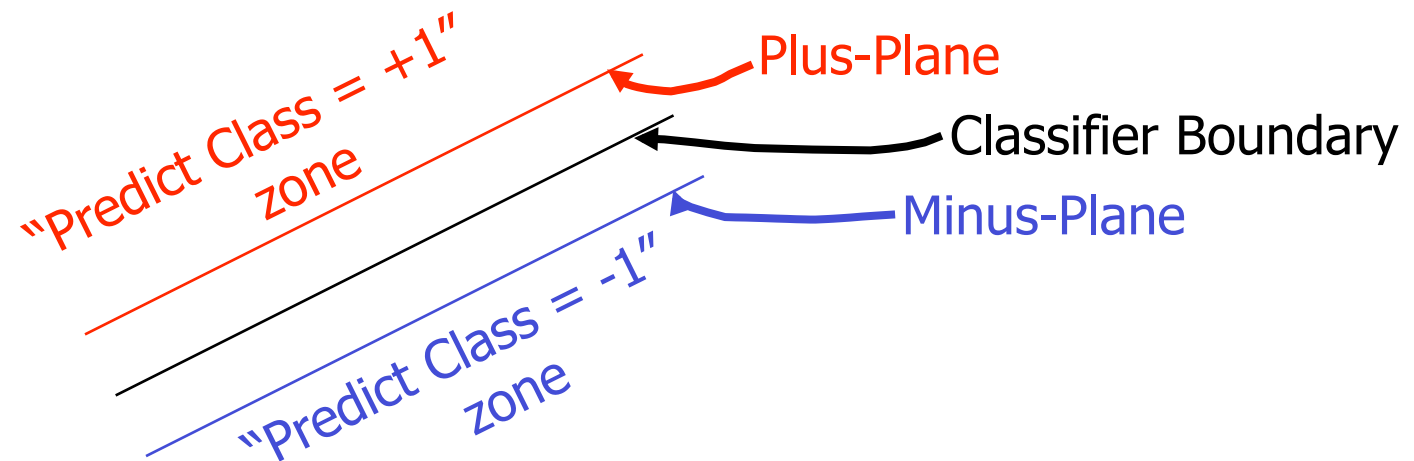
Linear SVM

# Why Maximum Margin?



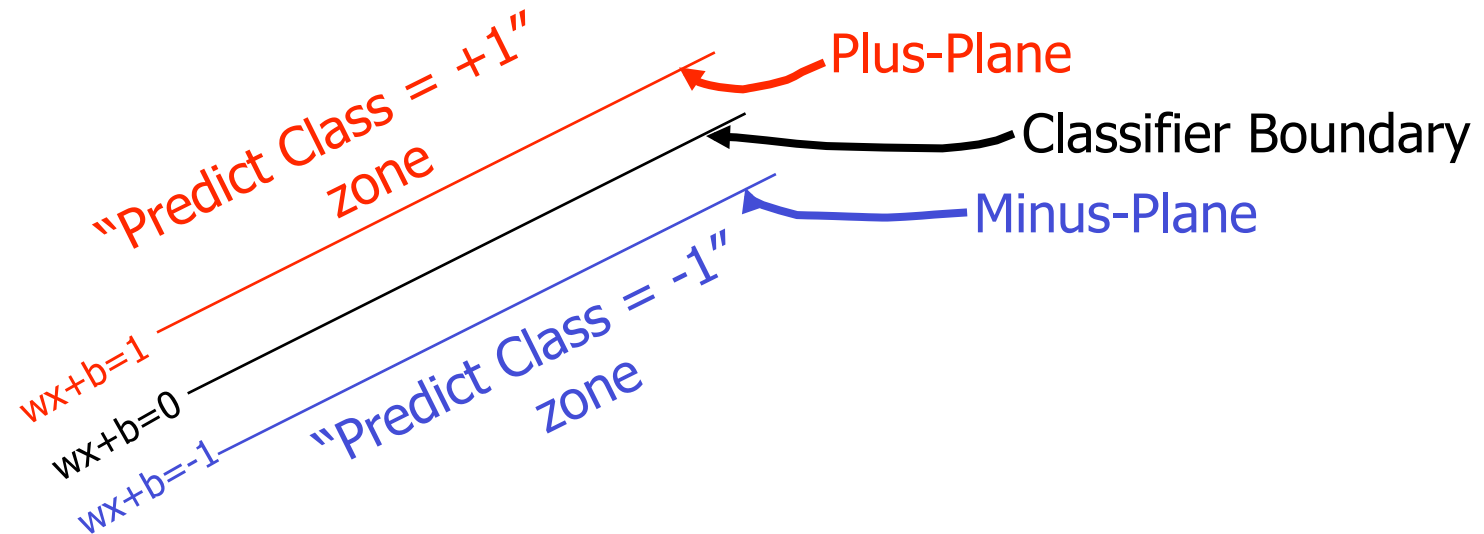
1. Intuitively this feels safest.
2. If we've made a small error in the location of the boundary (it's been jolted in its perpendicular direction) this gives us least chance of causing a misclassification.
3. LOOCV is easy since the model is immune to removal of any non-support-vector datapoints.
4. There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.
5. Empirically it works very very well.

# Specifying a line and margin



- How do we represent this mathematically?
- ...in  $m$  input dimensions?

# Specifying a line and margin



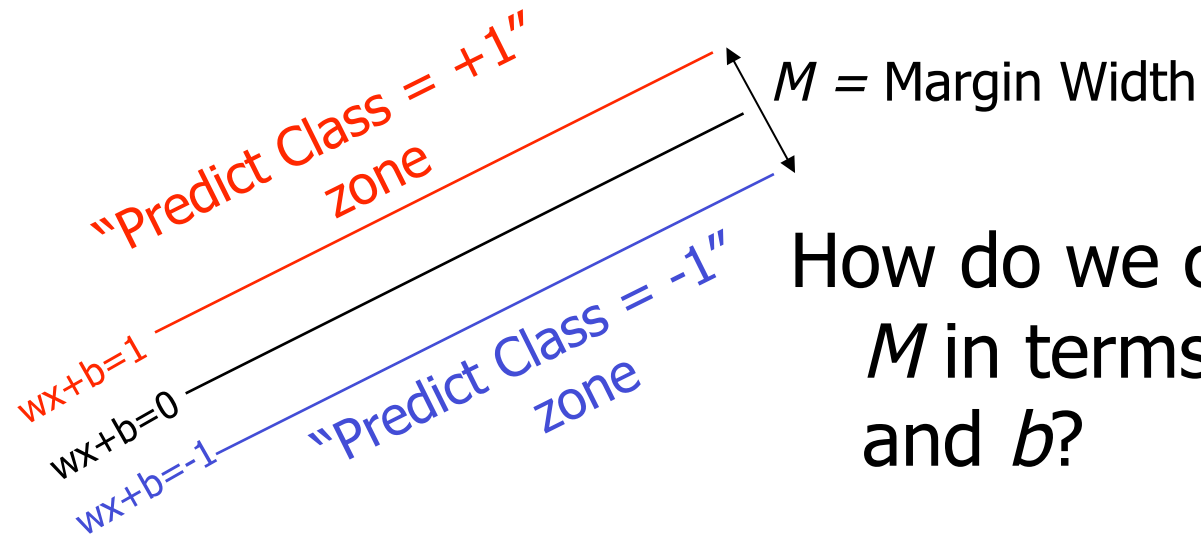
- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

Classify as..  $+1$  if  $\mathbf{w} \cdot \mathbf{x} + b \geq 1$

$-1$  if  $\mathbf{w} \cdot \mathbf{x} + b \leq -1$

Universe  
explodes if  $-1 < \mathbf{w} \cdot \mathbf{x} + b < 1$

# Computing the margin width

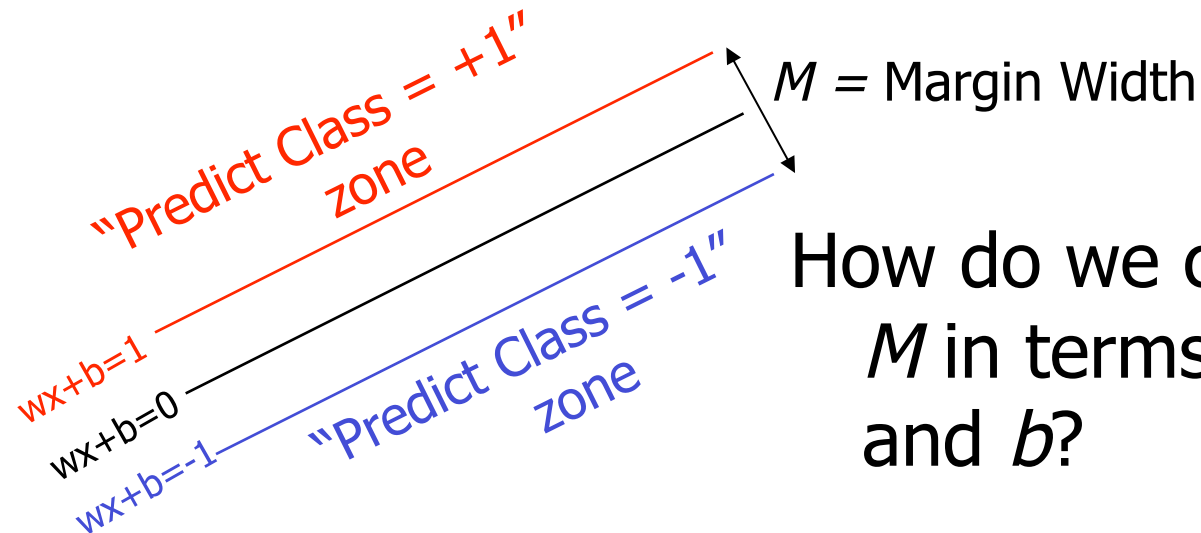


How do we compute  $M$  in terms of  $\mathbf{w}$  and  $b$ ?

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

**Claim:** The vector  $\mathbf{w}$  is perpendicular to the Plus Plane. **Why?**

# Computing the margin width



How do we compute  $M$  in terms of  $\mathbf{w}$  and  $b$ ?

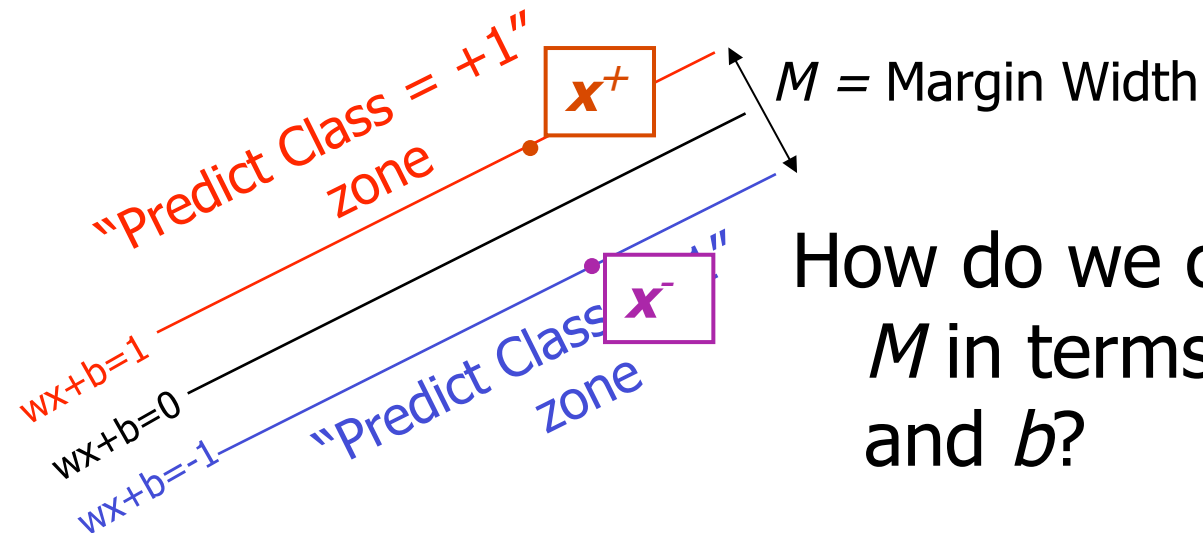
- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

**Claim:** The vector  $\mathbf{w}$  is perpendicular to the Plus Plane. **Why?**

Let  $\mathbf{u}$  and  $\mathbf{v}$  be two vectors on the Plus Plane. What is  $\mathbf{w} \cdot (\mathbf{u} - \mathbf{v})$ ?

And so of course the vector  $\mathbf{w}$  is also perpendicular to the Minus Plane

# Computing the margin width

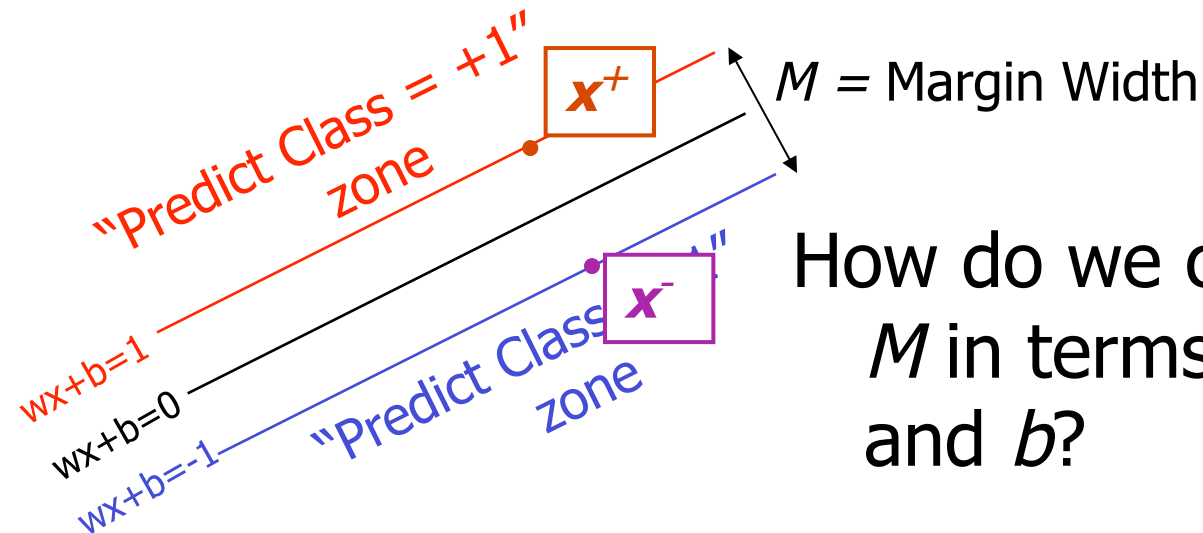


How do we compute  $M$  in terms of  $\mathbf{w}$  and  $b$ ?

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- The vector  $\mathbf{w}$  is perpendicular to the Plus Plane
- Let  $\mathbf{x}^-$  be any point on the minus plane
- Let  $\mathbf{x}^+$  be the closest plus-plane-point to  $\mathbf{x}^-$ .

Any location in  $\mathbb{R}^m$ : not necessarily a datapoint

# Computing the margin width

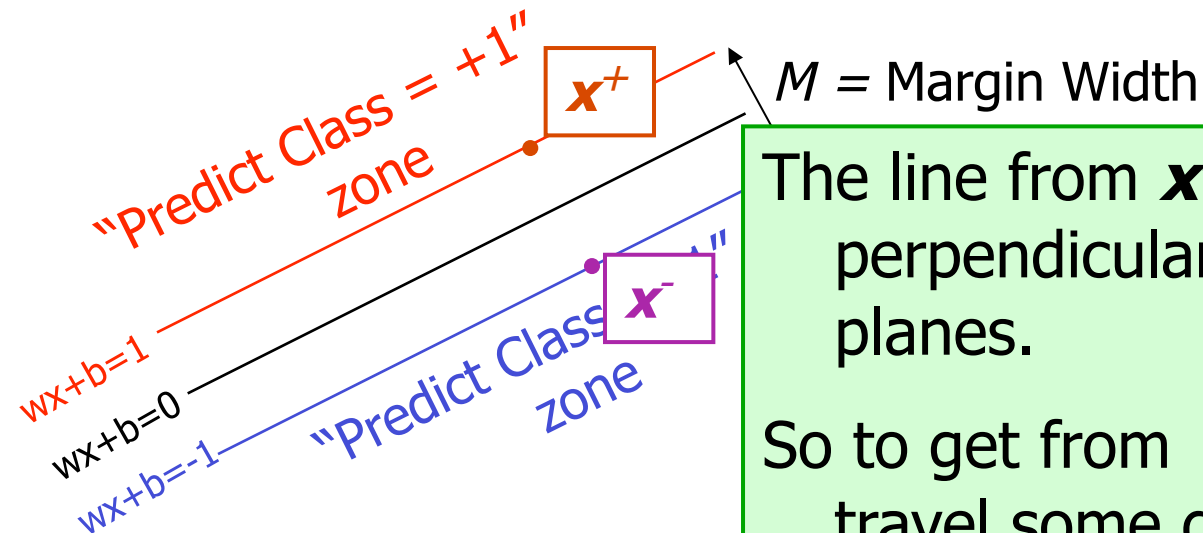


How do we compute  $M$  in terms of  $\mathbf{w}$  and  $b$ ?

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- The vector  $\mathbf{w}$  is perpendicular to the Plus Plane
- Let  $\mathbf{x}^-$  be any point on the minus plane
- Let  $\mathbf{x}^+$  be the closest plus-plane-point to  $\mathbf{x}^-$ .
- **Claim:**  $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$  for some value of  $\lambda$ . **Why?**



# Computing the margin width

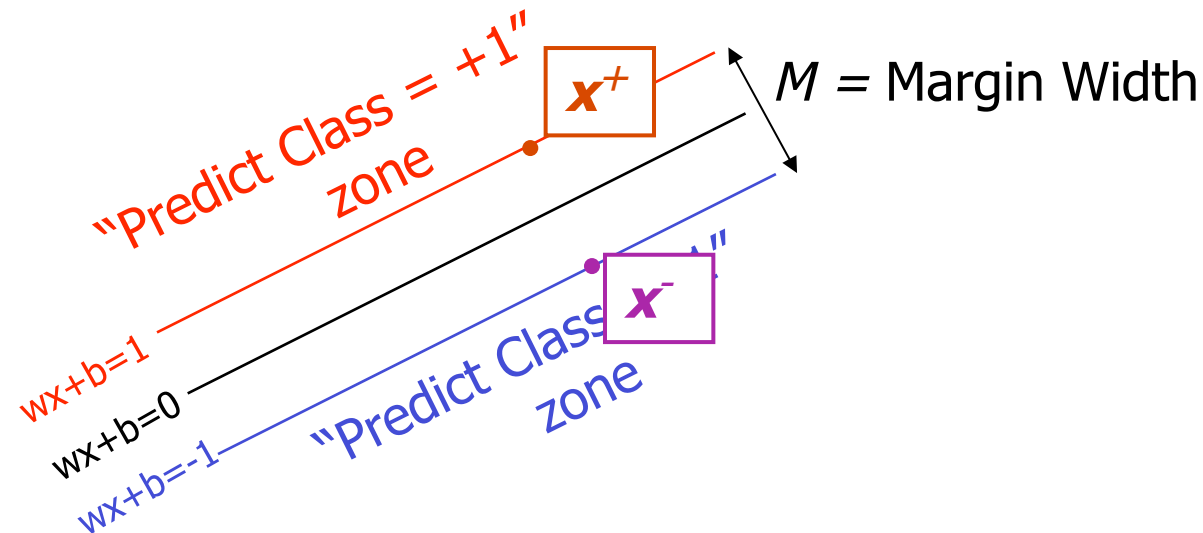


The line from  $\mathbf{x}^-$  to  $\mathbf{x}^+$  is perpendicular to the planes.

So to get from  $\mathbf{x}^-$  to  $\mathbf{x}^+$  travel some distance in direction  $\mathbf{w}$ .

- Plus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = 1 \}$
- Minus-plane =  $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- The vector  $\mathbf{w}$  is perpendicular to the Plus Plane
- Let  $\mathbf{x}^-$  be any point on the minus plane
- Let  $\mathbf{x}^+$  be the closest plus-plane-point to  $\mathbf{x}^-$ .
- **Claim:**  $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$  for some value of  $\lambda$ . **Why?**

# Computing the margin width

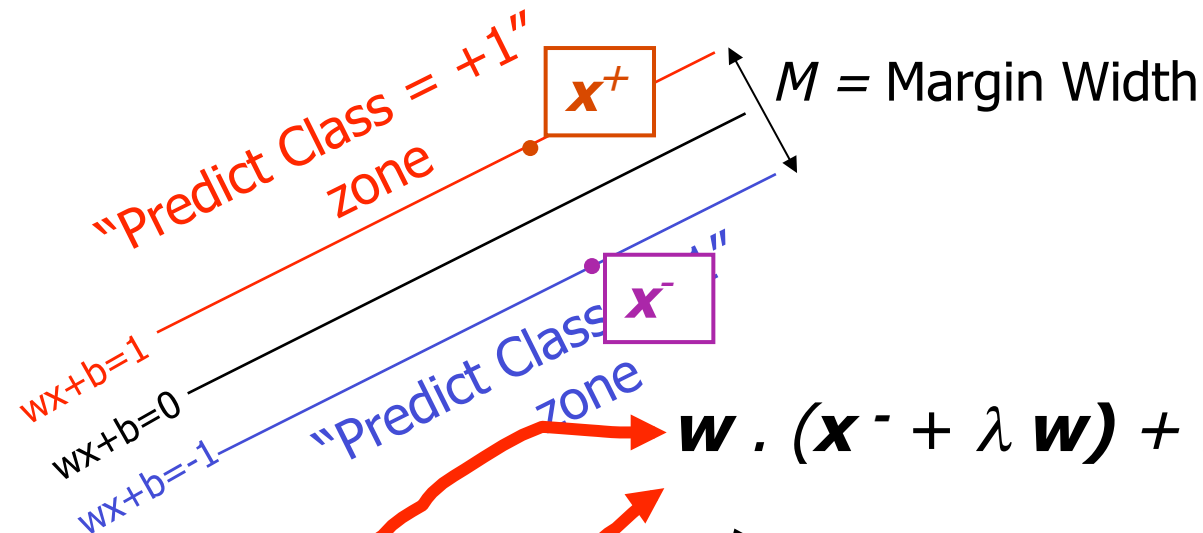


What we know:

- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$

It's now easy to get  $M$   
in terms of  $\mathbf{w}$  and  $b$

# Computing the margin width



What we know:

- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $x^+ = x^- + \lambda w$
- $|x^+ - x^-| = M$

It's now easy to get  $M$   
in terms of  $w$  and  $b$

$$w \cdot (x^- + \lambda w) + b = 1$$

$\Rightarrow$

$$w \cdot x^- + b + \lambda w \cdot w = 1$$

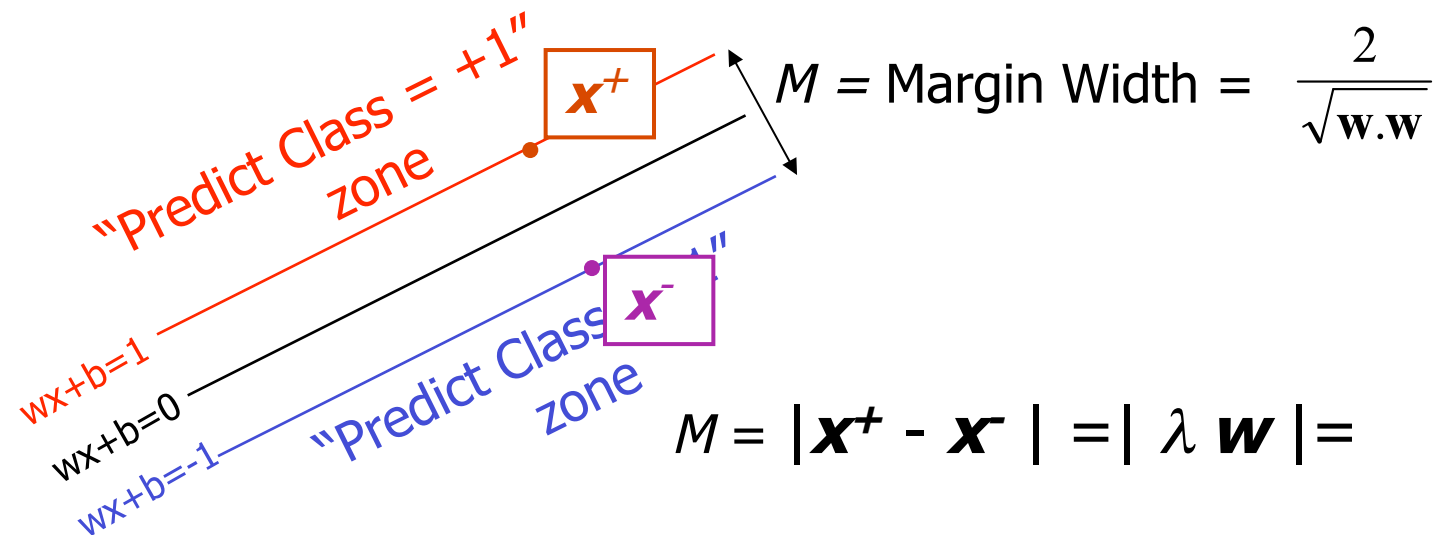
$\Rightarrow$

$$-1 + \lambda w \cdot w = 1$$

$\Rightarrow$

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

# Computing the margin width



$$M = |\mathbf{x}^+ - \mathbf{x}^-| = |\lambda \mathbf{w}| =$$

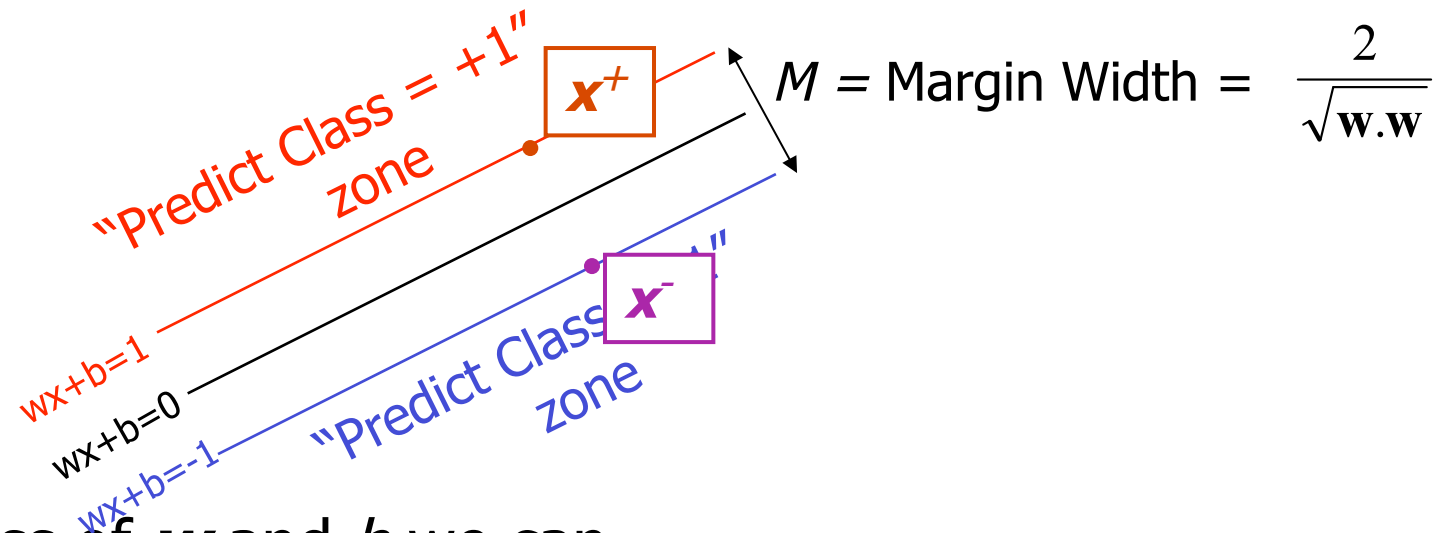
$$= \lambda |\mathbf{w}| = \lambda \sqrt{\mathbf{w} \cdot \mathbf{w}}$$

$$= \frac{2\sqrt{\mathbf{w} \cdot \mathbf{w}}}{\mathbf{w} \cdot \mathbf{w}} = \frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$$

What we know:

- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$
- $\mathbf{x}^+ = \mathbf{x}^- + \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$
- $\lambda = \frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$

# Learning the Maximum Margin Classifier



Given a guess of  $\mathbf{w}$  and  $b$  we can

- Compute whether all data points in the correct half-planes
- Compute the width of the margin

So now we just need to write a program to search the space of  $\mathbf{w}$ 's and  $b$ 's to find the widest margin that matches all the datapoints. *How?*

Gradient descent? Simulated Annealing? Matrix Inversion?  
EM? Newton's Method?

# Learning via Quadratic Programming

- QP is a well-studied class of optimization algorithms to maximize a quadratic function of some real-valued variables subject to linear constraints.

# Quadratic Programming

Find  $\arg \max_{\mathbf{u}} \quad c + \mathbf{d}^T \mathbf{u} + \frac{\mathbf{u}^T R \mathbf{u}}{2}$  ← Quadratic criterion

Subject to

$$\begin{aligned} a_{11}u_1 + a_{12}u_2 + \dots + a_{1m}u_m &\leq b_1 \\ a_{21}u_1 + a_{22}u_2 + \dots + a_{2m}u_m &\leq b_2 \\ &\vdots \\ a_{n1}u_1 + a_{n2}u_2 + \dots + a_{nm}u_m &\leq b_n \end{aligned}$$

}  $n$  additional linear inequality constraints

And subject to

$$\begin{aligned} a_{(n+1)1}u_1 + a_{(n+1)2}u_2 + \dots + a_{(n+1)m}u_m &= b_{(n+1)} \\ a_{(n+2)1}u_1 + a_{(n+2)2}u_2 + \dots + a_{(n+2)m}u_m &= b_{(n+2)} \\ &\vdots \\ a_{(n+e)1}u_1 + a_{(n+e)2}u_2 + \dots + a_{(n+e)m}u_m &= b_{(n+e)} \end{aligned}$$

}  $e$  additional linear equality constraints

# Quadratic Programming

Find  $\arg \max_{\mathbf{u}} \quad c + \mathbf{d}^T \mathbf{u} + \frac{\mathbf{u}^T \mathbf{R} \mathbf{u}}{2}$  ← Quadratic criterion

Subject to

There exist algorithms for finding such constrained quadratic optima much more efficiently and reliably than gradient ascent.

And subject to

(But they are very fiddly...you probably don't want to write one yourself)

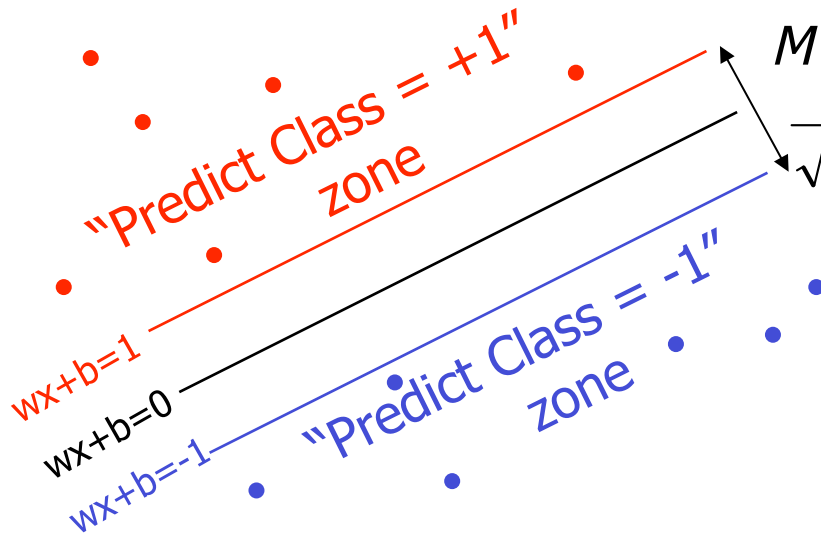
additional linear  
equality  
constraints

additional linear  
equality  
constraints

$$a_{(n+e)1}u_1 + a_{(n+e)2}u_2 + \dots + a_{(n+e)m}u_m = b_{(n+e)}$$



# Learning the Maximum Margin Classifier



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute whether all data points are in the correct half-planes
- Compute the margin width

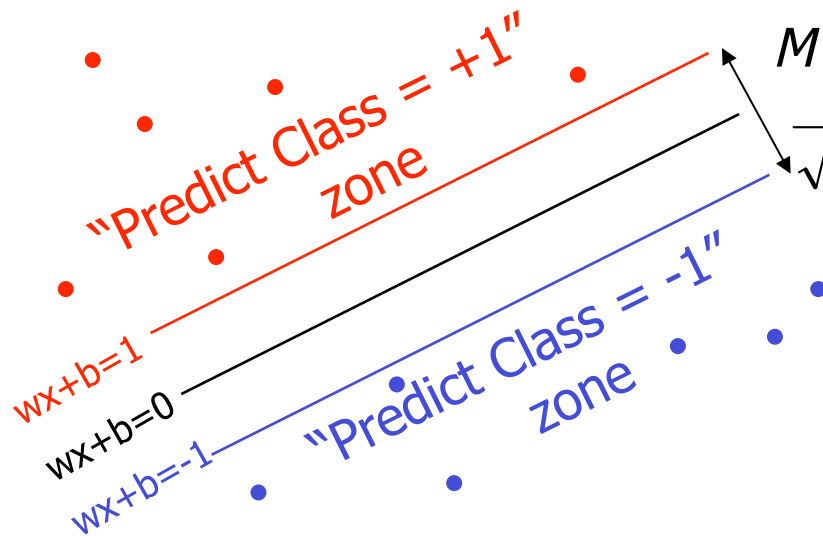
Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

What should our quadratic optimization criterion be?

How many constraints will we have?

What should they be?

# Learning the Maximum Margin Classifier



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute whether all data points are in the correct half-planes
- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

What should our quadratic optimization criterion be?

Minimize  $\mathbf{w} \cdot \mathbf{w}$

How many constraints will we have?  $R$

What should they be?

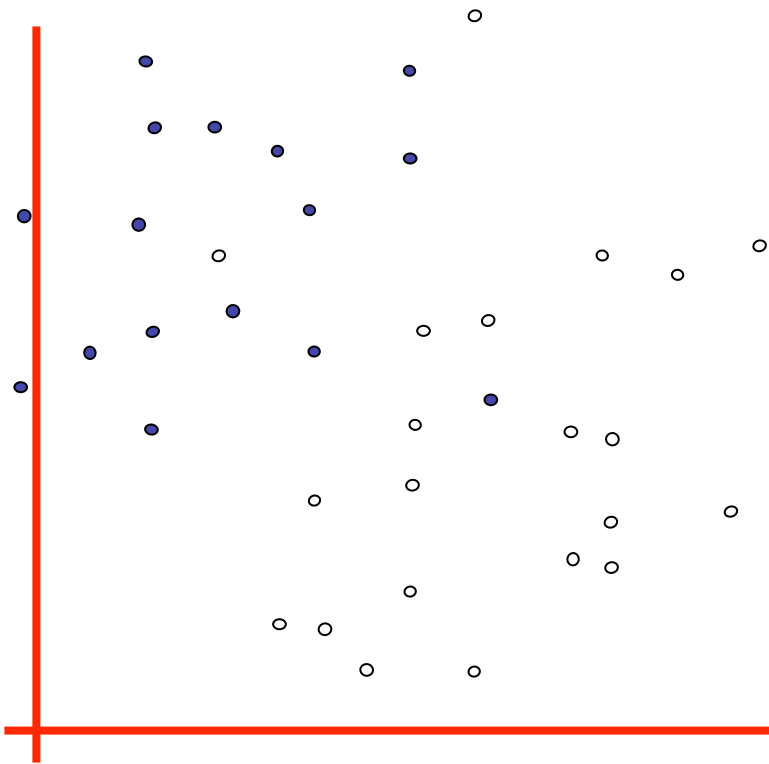
$$\mathbf{w} \cdot \mathbf{x}_k + b \geq 1 \text{ if } y_k = 1$$

$$\mathbf{w} \cdot \mathbf{x}_k + b \leq -1 \text{ if } y_k = -1$$

# Uh-oh!

This is going to be a problem!  
What should we do?

- denotes +1
- denotes -1

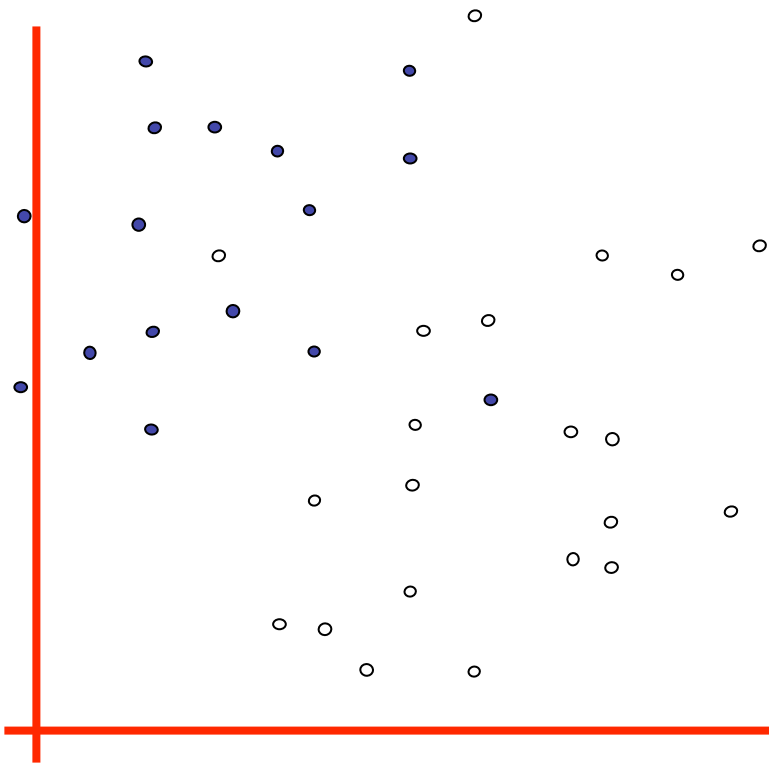


# Uh-oh!

This is going to be a problem!

What should we do?

- denotes +1
- denotes -1



Idea 1:

Find minimum  $\mathbf{w} \cdot \mathbf{w}$ , while minimizing number of training set errors.

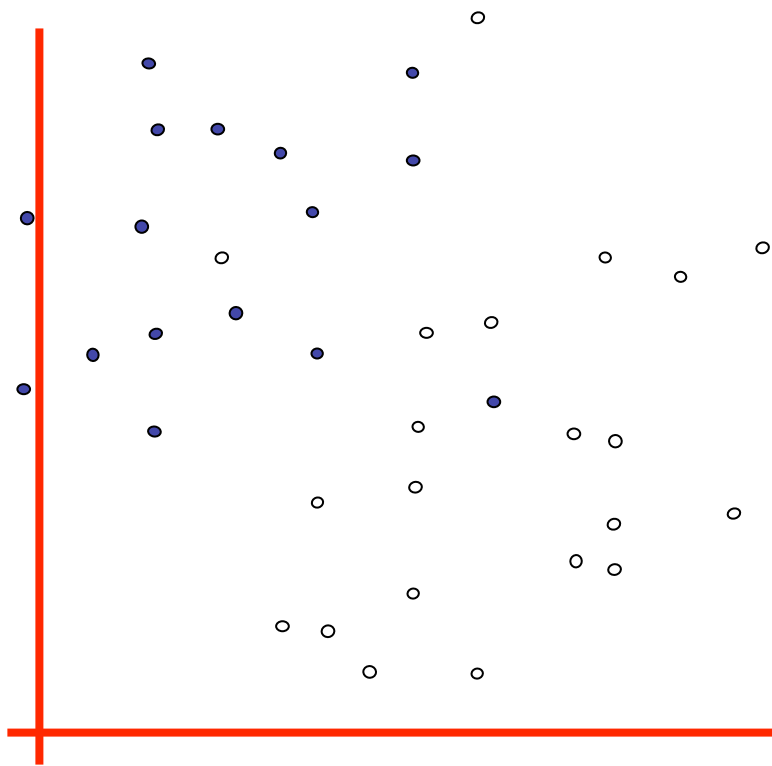
Problem: Two things to minimize makes for an ill-defined optimization

# Uh-oh!

This is going to be a problem!

What should we do?

- denotes +1
- denotes -1



Idea 1.1:

Minimize

$\mathbf{w} \cdot \mathbf{w} + C (\#train\ errors)$

Tradeoff parameter

There's a serious practical problem that's about to make us reject this approach. Can you guess what it is?

# Uh-oh!

This is going to be a problem!  
What should we do?

- denotes +1
- denotes -1

Idea 1.1:

Minimize

$\mathbf{w} \cdot \mathbf{w} + C (\#train\ errors)$

Tradeoff parameter

Can't be expressed as a Quadratic Programming problem.

Solving it may be too slow.

(Also, doesn't distinguish between disastrous errors and near misses)

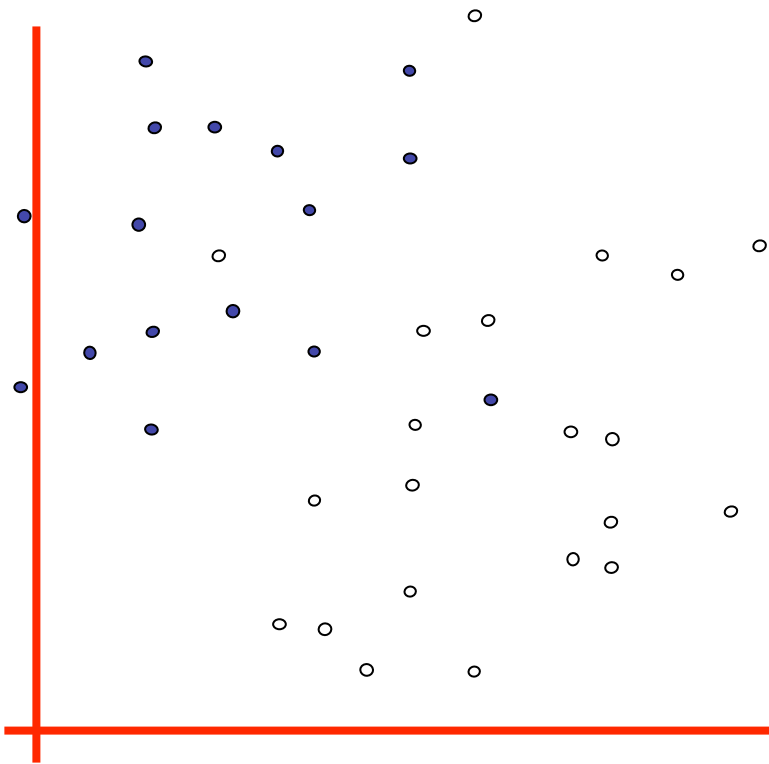
So... any other ideas?

# Uh-oh!

This is going to be a problem!

What should we do?

- denotes +1
- denotes -1

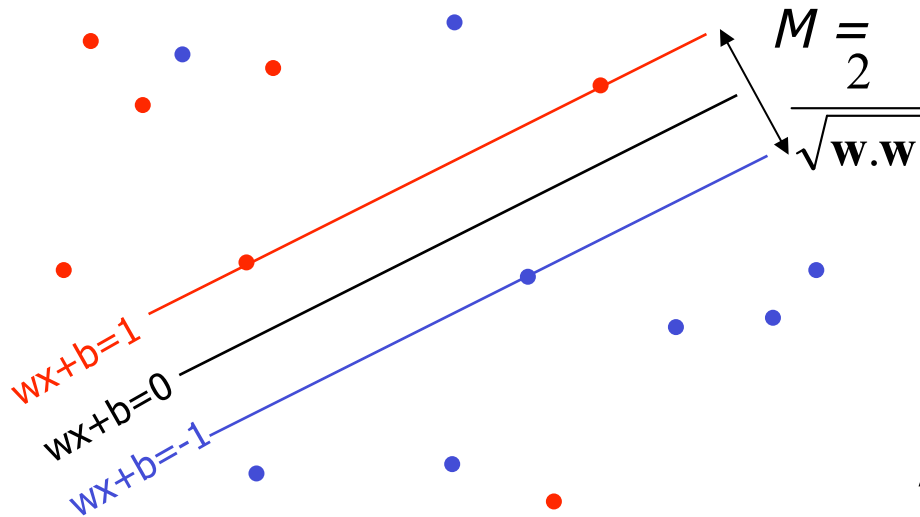


Idea 2.0:

Minimize

$\mathbf{W} \cdot \mathbf{W} + C$  (*distance of error points to their correct place*)

# Learning Maximum Margin with Noise



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute sum of distances of points to their correct zones

- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

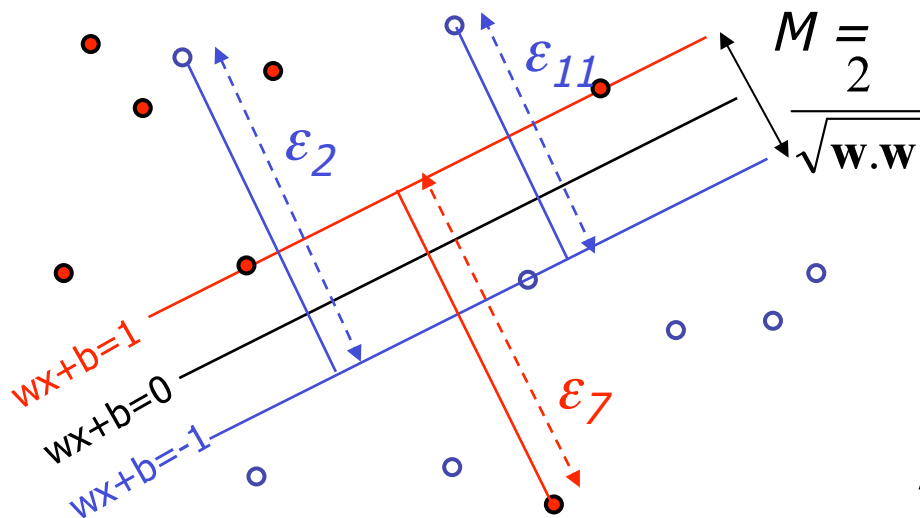
What should our quadratic optimization criterion be?

How many constraints will we have?

What should they be?



# Learning Maximum Margin with Noise



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute sum of distances of points to their correct zones

- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \varepsilon_k$$

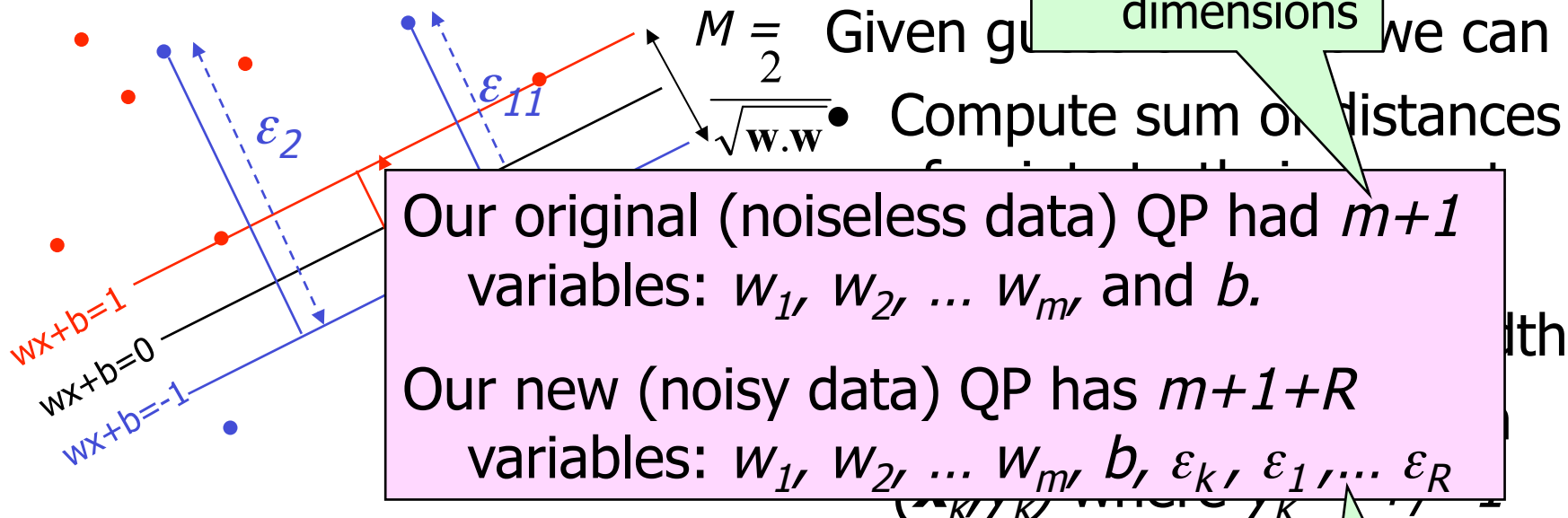
How many constraints will we have?  $R$

What should they be?

$$\mathbf{w} \cdot \mathbf{x}_k + b \geq 1 - \varepsilon_k \text{ if } y_k = 1$$

$$\mathbf{w} \cdot \mathbf{x}_k + b \leq -1 + \varepsilon_k \text{ if } y_k = -1$$

# Learning Maximum Margin with Noise



What should our quadratic optimization criterion be?

Minimize 
$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \varepsilon_k$$

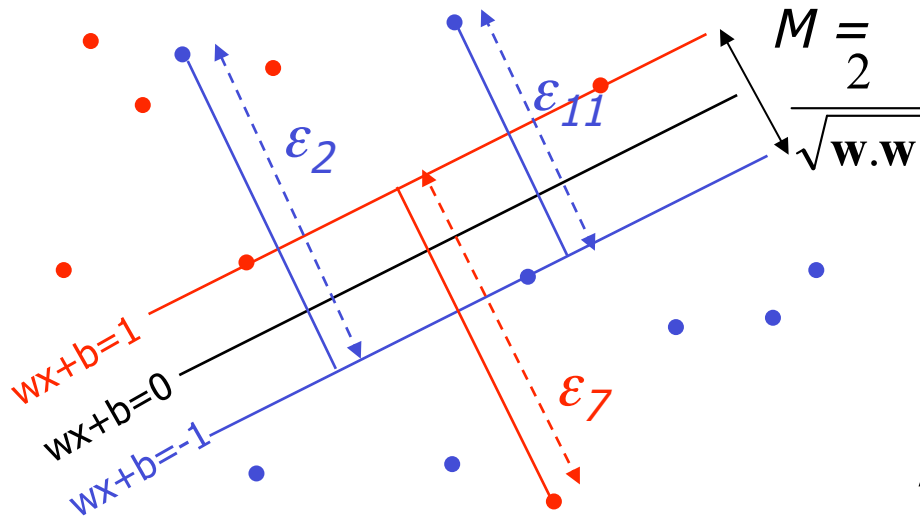
How many constraints do we have?  $R$

What should they be?

$\mathbf{w} \cdot \mathbf{x}_k + b \geq 1 - \varepsilon_k \text{ if } y_k = 1$

$\mathbf{w} \cdot \mathbf{x}_k + b \leq -1 + \varepsilon_k \text{ if } y_k = -1$

# Learning Maximum Margin with Noise



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute sum of distances of points to their correct zones

- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

What should our quadratic optimization criterion be?

Minimize 
$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \epsilon_k$$

How many constraints will we have?  $R$

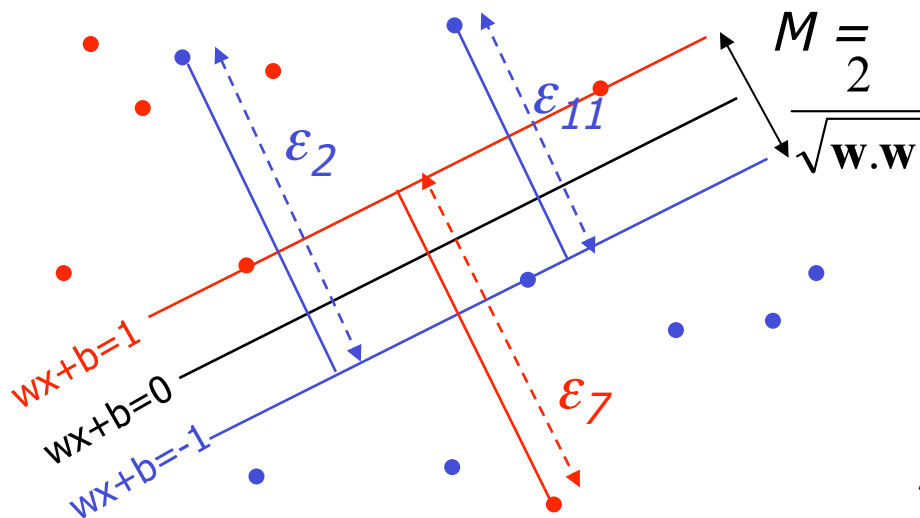
What should they be?

$$\mathbf{w} \cdot \mathbf{x}_k + b \geq 1 - \epsilon_k \text{ if } y_k = 1$$

$$\mathbf{w} \cdot \mathbf{x}_k + b \leq -1 + \epsilon_k \text{ if } y_k = -1$$

There's a bug in this QP. Can you spot it?

# Learning Maximum Margin with Noise



Given guess of  $\mathbf{w}$ ,  $b$  we can

- Compute sum of distances of points to their correct zones

- Compute the margin width

Assume  $R$  datapoints, each  $(\mathbf{x}_k, y_k)$  where  $y_k = +/- 1$

What should our quadratic optimization criterion be?

Minimize 
$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \epsilon_k$$

How many constraints will we have?  $2R$

What should they be?

$$\mathbf{w} \cdot \mathbf{x}_k + b \geq 1 - \epsilon_k \text{ if } y_k = 1$$

$$\mathbf{w} \cdot \mathbf{x}_k + b \leq -1 + \epsilon_k \text{ if } y_k = -1$$

$$\epsilon_k \geq 0 \text{ for all } k$$