

Support Vector Machines

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.

Andrew W. Moore
Professor
School of Computer Science
Carnegie Mellon University

www.cs.cmu.edu/~awm

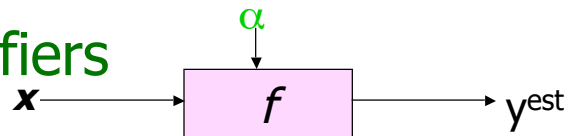
awm@cs.cmu.edu

412-268-7599

Copyright © 2001, 2003, Andrew W. Moore

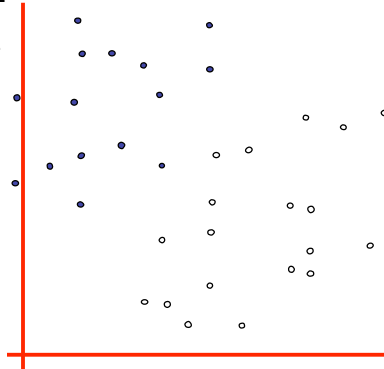
Nov 23rd, 2001

Linear Classifiers



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

- denotes +1
- denotes -1

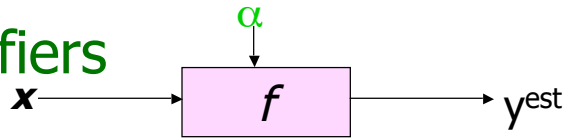


How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

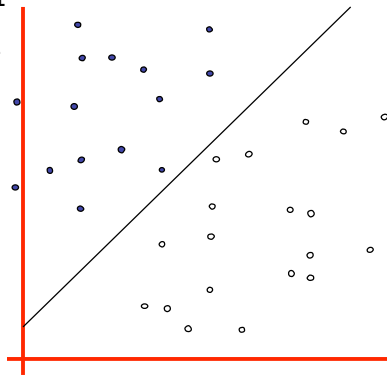
Support Vector Machines: Slide 2

Linear Classifiers



- denotes +1
- denotes -1

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

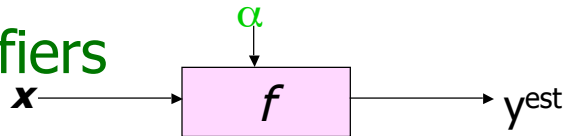


How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

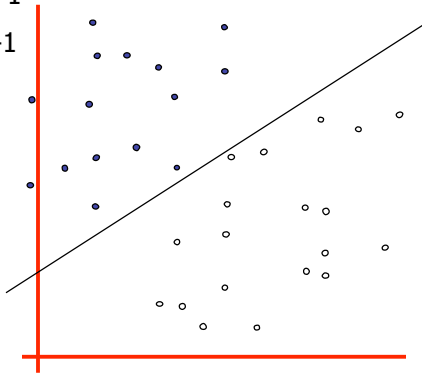
Support Vector Machines: Slide 3

Linear Classifiers



- denotes +1
- denotes -1

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

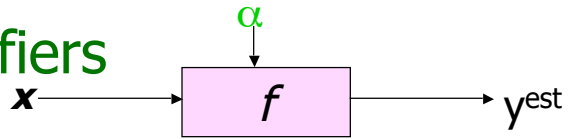


How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

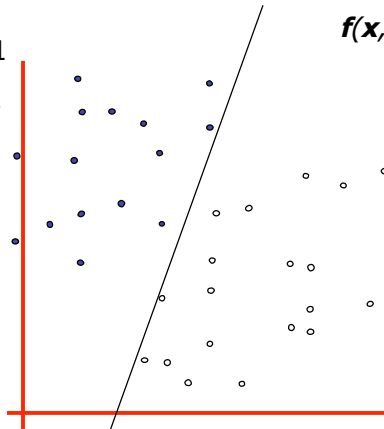
Support Vector Machines: Slide 4

Linear Classifiers



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

- denotes +1
- denotes -1

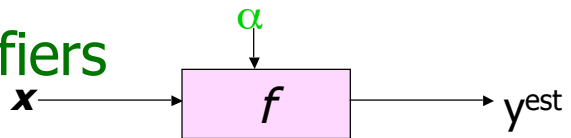


How would you classify this data?

Copyright © 2001, 2003, Andrew W. Moore

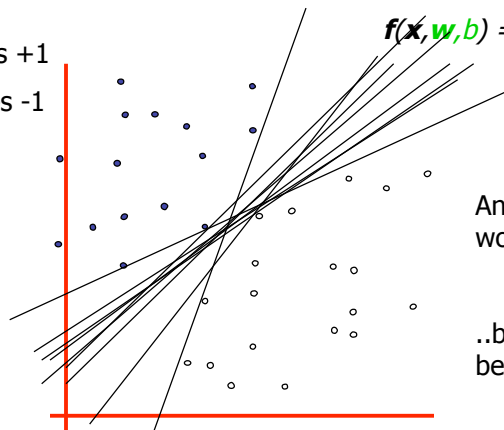
Support Vector Machines: Slide 5

Linear Classifiers



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

- denotes +1
- denotes -1



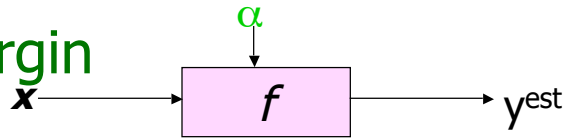
Any of these would be fine..

..but which is best?

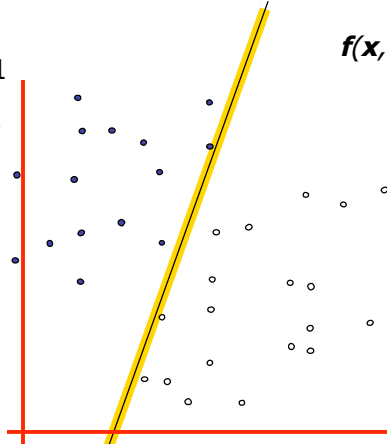
Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 6

Classifier Margin



- denotes +1
- denotes -1



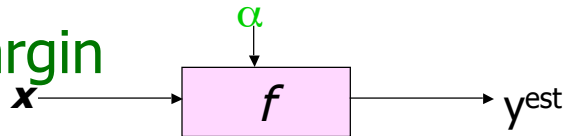
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

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

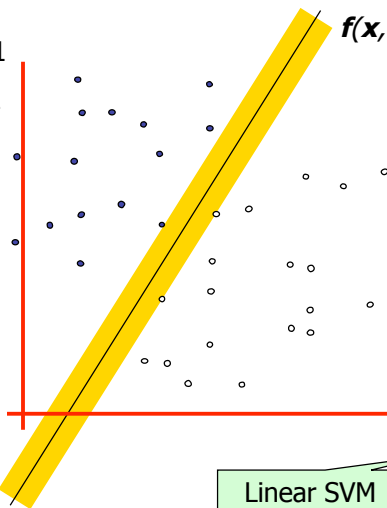
Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 7

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)

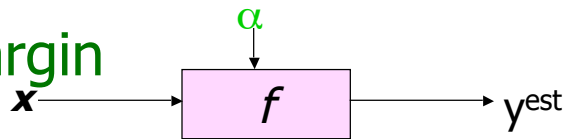
Note: **SMO** (Sequential Minimal Optimization) is one algorithm for computing \mathbf{w} and b .

Linear SVM

Copyright © 2001, 2003, Andrew W. Moore

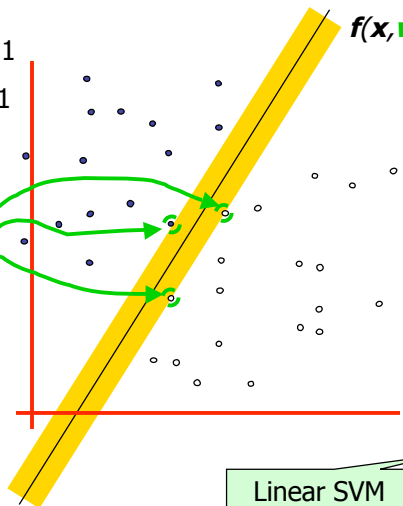
Support Vector Machines: Slide 8

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.

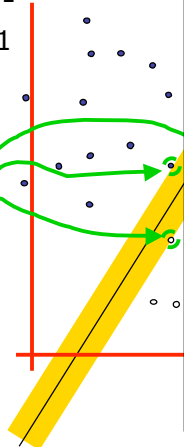
This is the simplest kind of SVM (Called an LSVM)

Linear SVM

Why Maximum Margin?

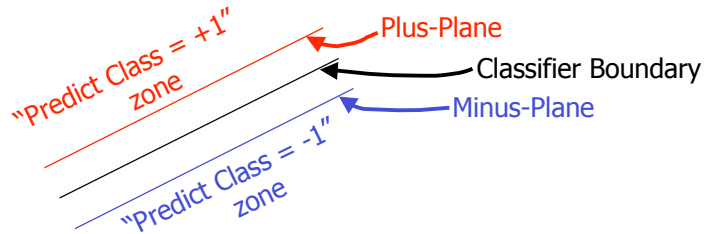
- denotes +1
- denotes -1

Support Vectors are those datapoints that the margin pushes up against



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

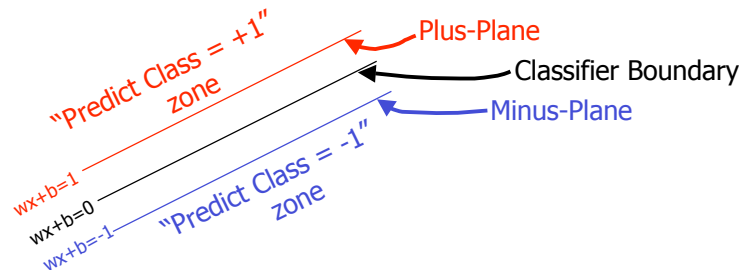


- How do we represent this mathematically?
- ...in m input dimensions?
- ...so that we can maximize the margin?

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 11

Specifying a line and margin



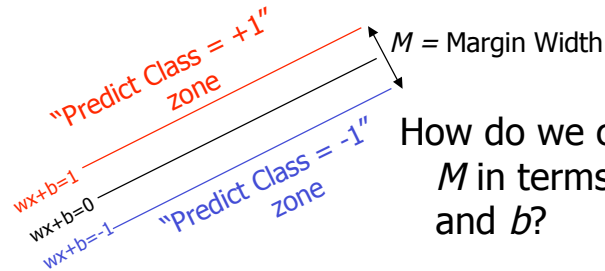
- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$

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

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 12

Computing the margin width

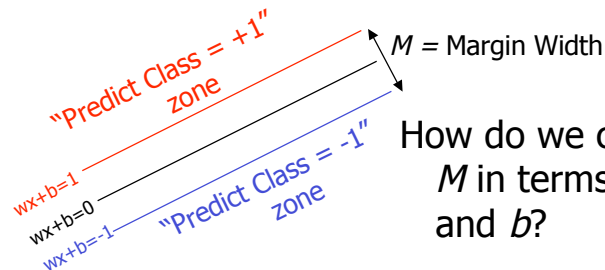


How do we compute M in terms of \mathbf{w} and b ?

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

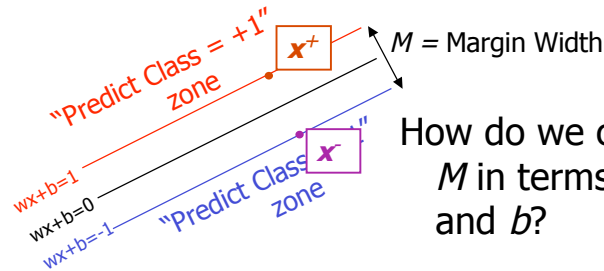
Claim: The vector \mathbf{w} is perpendicular to the Plus Plane. **Why?**

- Definitions: "vector" == "point"
- \mathbf{x}_1 perpendicular to \mathbf{x}_2 iff $\mathbf{x}_1 \cdot \mathbf{x}_2 == 0$

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 ?

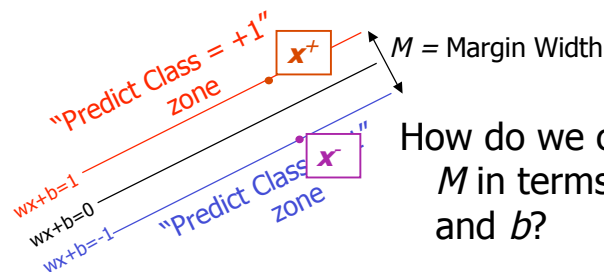
- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = +1 \}$
- Minus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = -1 \}$
- The vector \mathbf{w} is perpendicular to the Plus Plane
- Let \mathbf{x}^* be any point on the minus plane
- Let \mathbf{x}^\dagger be the closest plus-plane-point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

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

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 15

Computing the margin width



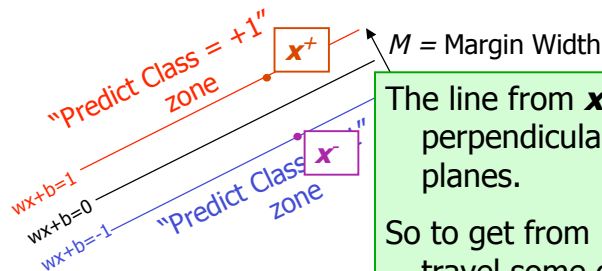
How do we compute M in terms of \mathbf{w} and b ?

- 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}^\dagger be the closest plus-plane-point to \mathbf{x}^* .
- **Claim:** $\mathbf{x}^\dagger = \mathbf{x}^* + \lambda \mathbf{w}$ for some value of λ . **Why?**
 - **Note:** λ is scalar and positive.

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 16

Computing the margin width



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

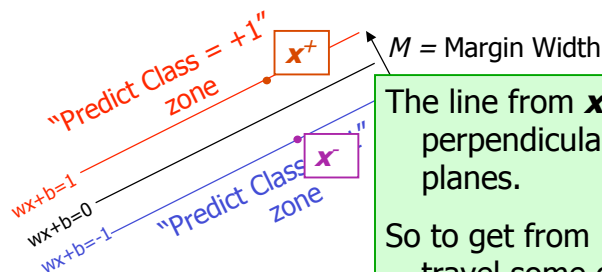
So to get from \mathbf{x} to \mathbf{x}^+ travel some distance in the direction of \mathbf{w} .

- 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 λ . **Why?**
 - **Note:** λ is scalar and positive.

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 17

Computing the margin width



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

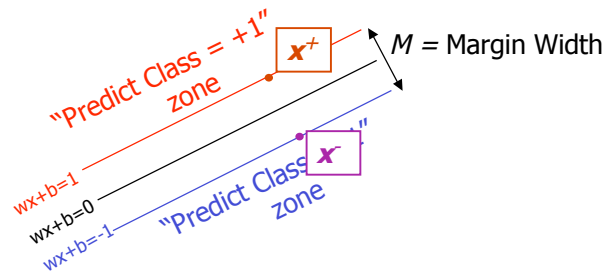
So to get from \mathbf{x} to \mathbf{x}^+ travel some distance in the direction of \mathbf{w} .

- Plus-plane = $\{ \mathbf{x} : \mathbf{w} \cdot \mathbf{x} + b = 1 \}$
- So to now we know that $\mathbf{x}^+ - \mathbf{x} = \lambda \mathbf{w}$.
- 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 λ . **Why?**
 - **Note:** λ is scalar and positive.

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 18

Computing the margin width



What we know:

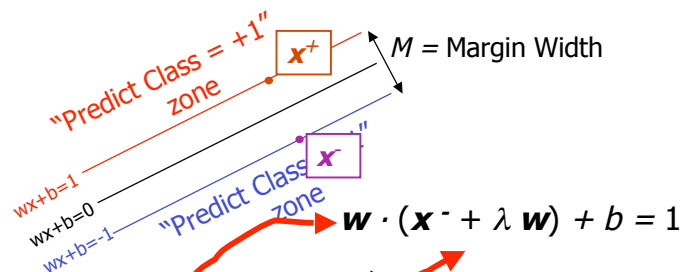
- $\mathbf{w} \cdot \mathbf{x}^+ + b = +1$
- $\mathbf{w} \cdot \mathbf{x}^- + b = -1$
- $\mathbf{x}^+ - \mathbf{x}^- = \lambda \mathbf{w}$
- $|\mathbf{x}^+ - \mathbf{x}^-| = M$

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

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 19

Computing the margin width



What we know:

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

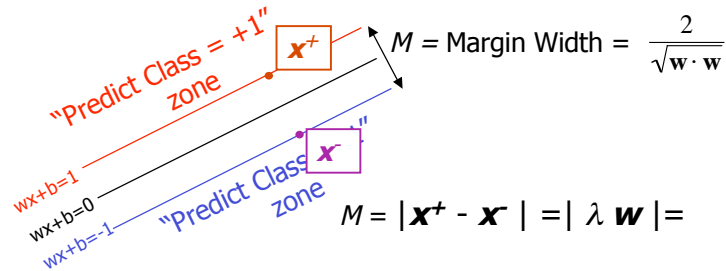
It's now easy to get M
in terms of \mathbf{w} and b

$$\begin{aligned} & \Rightarrow \mathbf{w} \cdot (\mathbf{x}^- + \lambda \mathbf{w}) + b = 1 \\ & (\mathbf{w} \cdot \mathbf{x}^- + b) + \lambda \mathbf{w} \cdot \mathbf{w} = 1 \\ & \Rightarrow \\ & -1 + \lambda \mathbf{w} \cdot \mathbf{w} = 1 \\ & \Rightarrow \lambda = \frac{2}{\mathbf{w} \cdot \mathbf{w}} \end{aligned}$$

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 20

Computing the margin width



What we know:

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

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

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

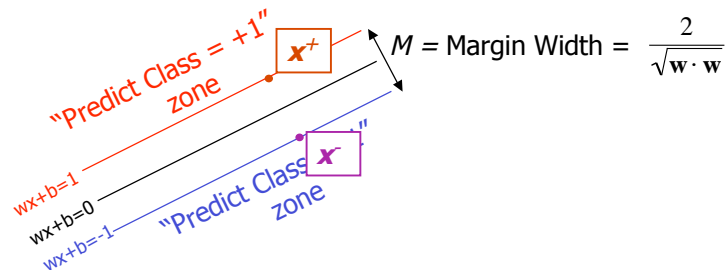
Yay! Just maximize $\frac{2}{\sqrt{\mathbf{w} \cdot \mathbf{w}}}$

Wait...OMG the data! I forgot!

Support Vector Machines: Slide 21

Copyright © 2001, 2003, Andrew W. Moore

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?

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 22

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.
- It will solve our problem for us!
- It doesn't matter how it works!
- Popular ML approach:
 - Describe your learning problem as optimization...
 - ...and give it to somebody else to solve!

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 23

Quadratic Programming - Don't Panic

Find $\underset{\mathbf{w}}{\operatorname{argmin}} \quad c + \mathbf{d}^T \mathbf{w} + \frac{\mathbf{w}^T \mathbf{K} \mathbf{w}}{2}$ ← Quadratic criterion
 Note $\mathbf{w} \cdot \mathbf{x} = \mathbf{w}^T \mathbf{x}$

Subject to

$$\left. \begin{aligned} a_{11}w_1 + a_{12}w_2 + \dots + a_{1m}w_m &\leq b_1 \\ a_{21}w_1 + a_{22}w_2 + \dots + a_{2m}w_m &\leq b_2 \\ &\vdots \\ a_{n1}w_1 + a_{n2}w_2 + \dots + a_{nm}w_m &\leq b_n \end{aligned} \right\} \begin{array}{l} n \text{ additional linear} \\ \text{inequality} \\ \text{constraints} \end{array}$$

And subject to

$$\left. \begin{aligned} a_{(n+1)1}w_1 + a_{(n+1)2}w_2 + \dots + a_{(n+1)m}w_m &= b_{(n+1)} \\ a_{(n+2)1}w_1 + a_{(n+2)2}w_2 + \dots + a_{(n+2)m}w_m &= b_{(n+2)} \\ &\vdots \\ a_{(n+e)1}w_1 + a_{(n+e)2}w_2 + \dots + a_{(n+e)m}w_m &= b_{(n+e)} \end{aligned} \right\} \begin{array}{l} e \text{ additional linear} \\ \text{equality} \\ \text{constraints} \end{array}$$

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 24

Quadratic Programming

Find $\arg \min_{\mathbf{w}} c + \mathbf{d}^T \mathbf{w} + \frac{\mathbf{w}^T \mathbf{K} \mathbf{w}}{2}$ ← Quadratic criterion

Subject to

And subject to

$$a_{(n+e)1} w_1 + a_{(n+e)2} w_2 + \dots + a_{(n+e)m} w_m = b_{(n+e)}$$

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

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

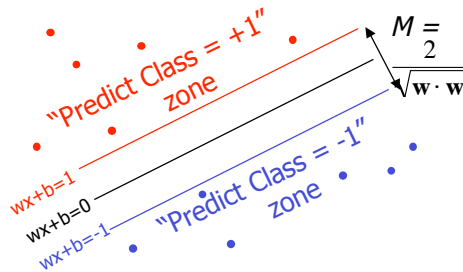
Additional linear equality constraints

e additional linear equality constraints

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 25

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?

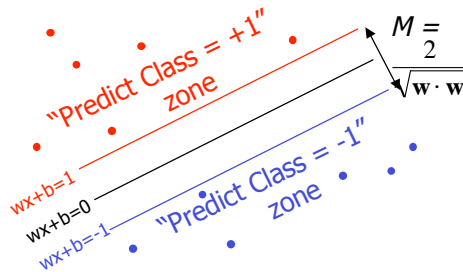
How many constraints will we have?

What should they be?

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 26

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?

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

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

Copyright © 2001, 2003, Andrew W. Moore

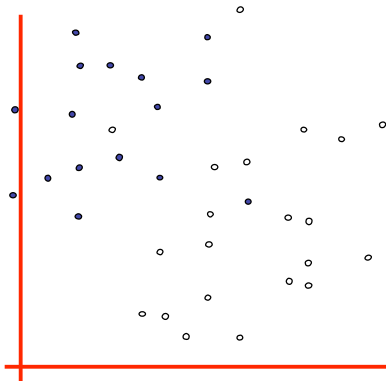
Support Vector Machines: Slide 27

Uh-oh!

This is going to be a problem!

What should we do?

- denotes +1
- denotes -1



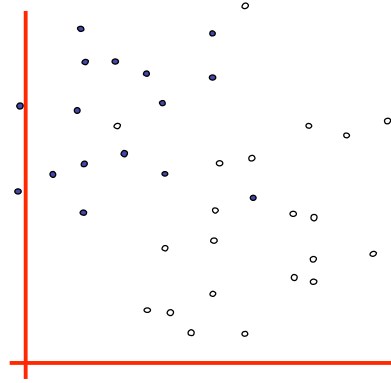
Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 28

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

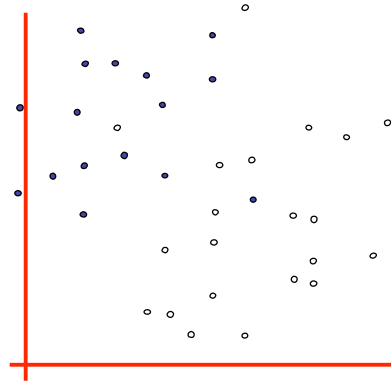
Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 29

Uh-oh!

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

- denotes +1
- denotes -1



Idea 1.1:

Minimize

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

Tradeoff parameter

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

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 30

Uh-oh!

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

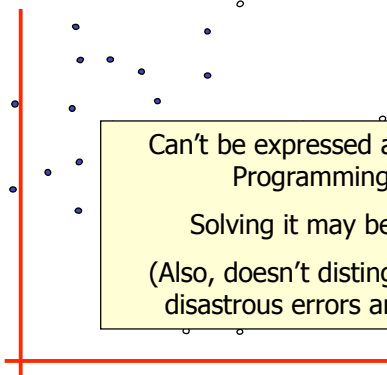
- denotes +1
- denotes -1

Idea 1.1:

Minimize

$$w \cdot w + C (\#train\ errors)$$

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!

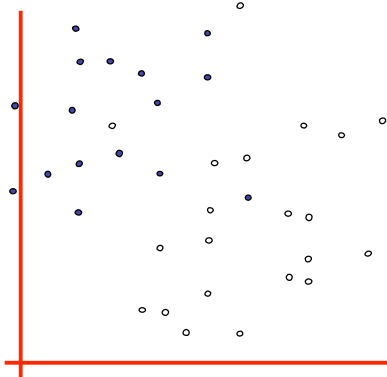
This is going to be a problem!
What should we do?

- denotes +1
- denotes -1

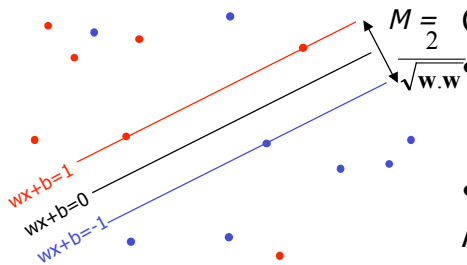
Idea 2.0:

Minimize

$$w \cdot w + C (\text{distance from incorrectly labeled points to their correct place})$$



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?

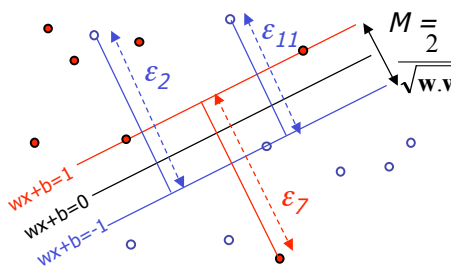
How many constraints will we have?

What should they be?

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 33

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? R

What should they be?

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

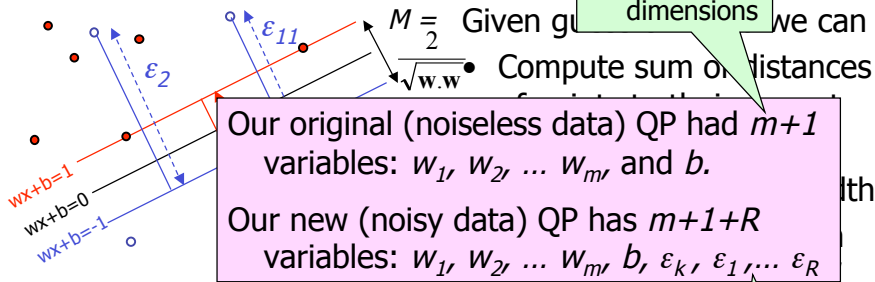
$$\mathbf{w} \cdot \mathbf{x}_k + b \geq (1 - \varepsilon_k) \text{ if } y_k = 1$$

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

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 34

Learning Maximum Margin with Noise



What should our quadratic optimization criterion be?

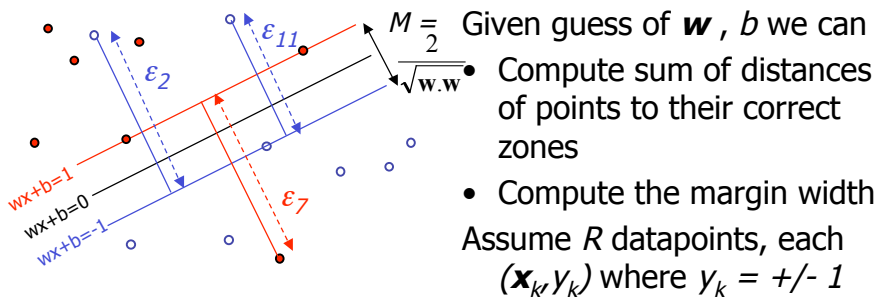
Minimize $\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \varepsilon_k$

How many constraints will we have? R

What should they be?

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

Learning Maximum Margin with Noise



What should our quadratic optimization criterion be?

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

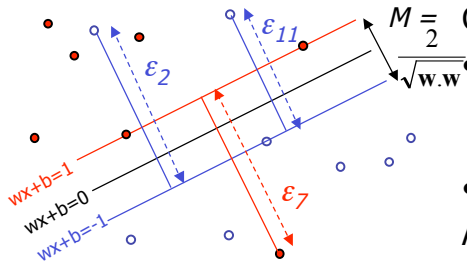
How many constraints will we have? R

What should they be?

$\mathbf{w} \cdot \mathbf{x}_k + b \geq (1 - \varepsilon_k)$ if $y_k = 1$
 $\mathbf{w} \cdot \mathbf{x}_k + b \leq (-1 + \varepsilon_k)$ 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?

How many constraints will we have? $2R$

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

What should they be?

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

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

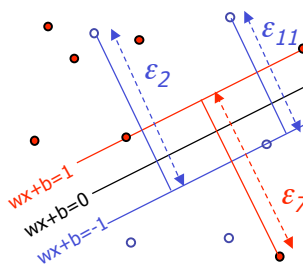
Called "slack variables"

$\varepsilon_k \geq 0$ for all k

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 37

Learning Maximum Margin with Noise



Big C means "Fit the training data as much as possible!" (at the expense of maximizing margin)

Small C means "Maximize the margin as much as possible!" (at the expense of fitting the training data)

, b we can
of distances
ir correct
margin width
nts, each
 $y_k = +/- 1$

What should our quadratic optimization criterion be?

How many constraints will we have? $2R$

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

What should they be?

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

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

Called "slack variables"

$\varepsilon_k \geq 0$ for all k

Copyright © 2001, 2003, Andrew W. Moore

Support Vector Machines: Slide 38

What do we have?

- Method for learning a maximum-margin linear classifier when the data are
 - “Linearly separable” - i.e. there is a line that gets 0 training error
 - Not linearly separable - there is no such line.
- If not linearly separable, we make a trade-off between maximizing margin and minimizing “stuff-is-on-the-wrong-side-ness”
- Still, our output for a given \mathbf{x} is

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