

# Support Vector Machines (SVMs)

Lyle Ungar

## Learning objectives

Review decision boundaries

Hinge loss

Margin

SVM – primal and dual

# Representing Lines

- How do we represent a line?

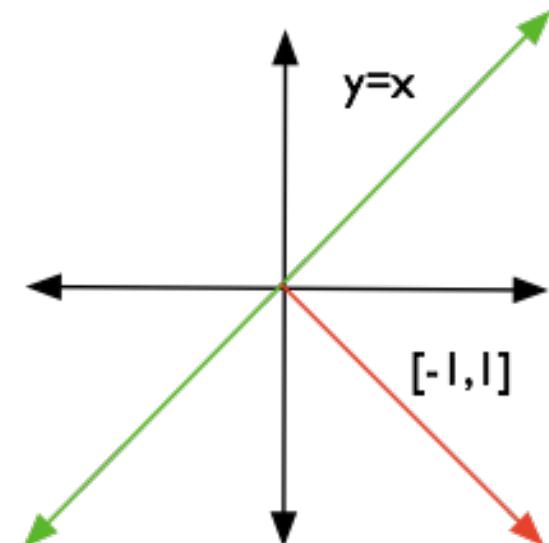
$$y = x$$

$$0 = x - y$$

$$0 = [1, -1] \begin{bmatrix} x \\ y \end{bmatrix}$$

- In general, a hyperplane is defined by

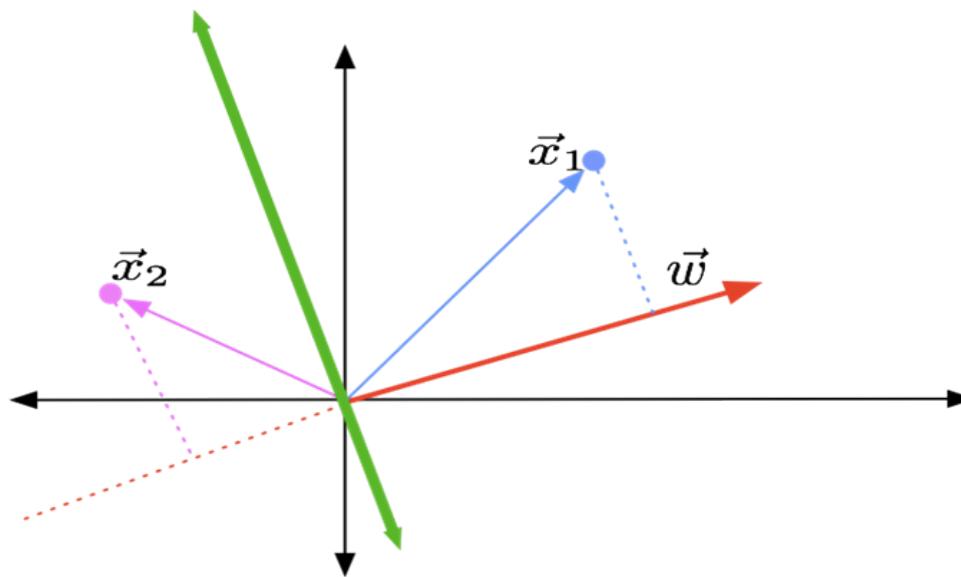
$$0 = \vec{w} \cdot \vec{x}$$



The red vector ( $w$ ) defines the green plane that is orthogonal to it.

Why bother with this weird representation?

# Projections



$(\vec{w} \cdot \vec{x})\vec{w}$  is the projection of  $\vec{x}$  onto  $\vec{w}$

alternate intuition: recall the dot product of two vectors is simply the product of their lengths and the cosine of the angle between them

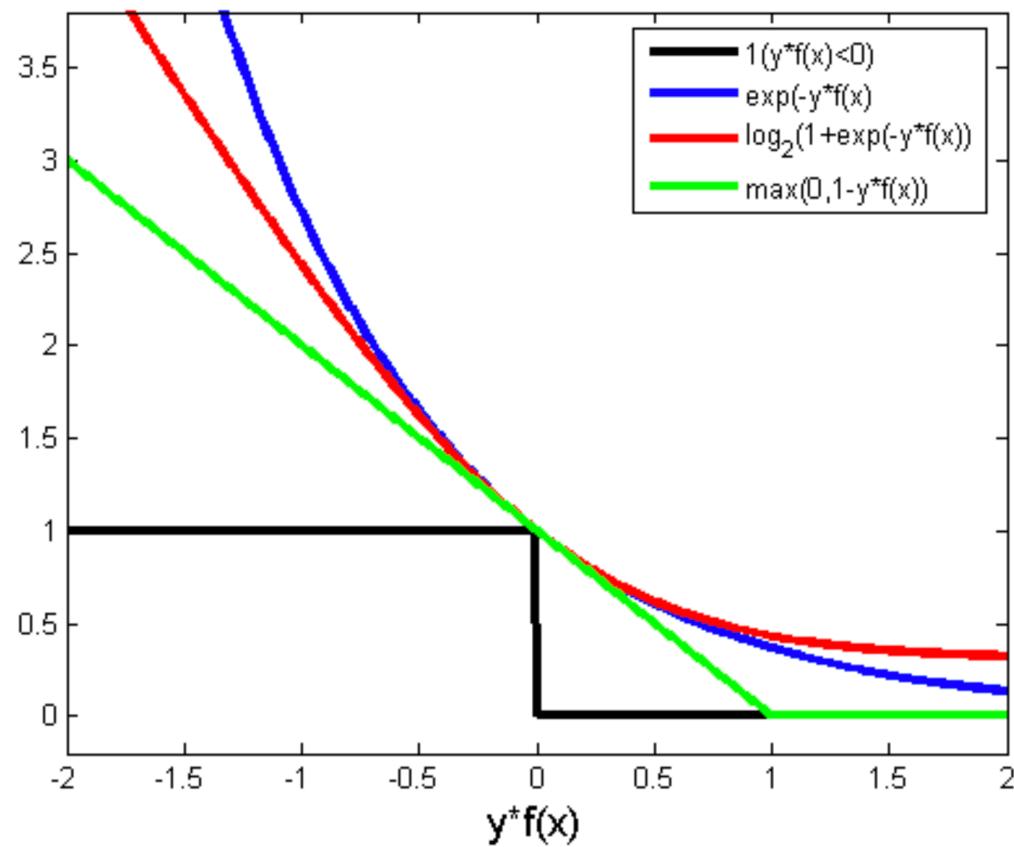
# Now classification is easy!

- ◆ Input:  $x$
- ◆ Model:  $w$
- ◆ Score:  $w^T x$
- ◆ Prediction:  $\text{sgn}(w^T x)$
- ◆ But how do we learn  $w$ ?

# Support Vector Machines (SVMs)

- ◆ Minimize hinge loss
  - With regularization
- ◆ "Large margin" methods

# Loss functions for classification



Hinge loss

# SVM: Hinge loss, ridge penalty

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

$$\min_{\mathbf{w}, b, \xi \geq 0} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_i \xi_i$$

$$\xi_i = \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b))$$

0 if score is right by 1 or more  
(hinge loss)

# Support Vector Machines

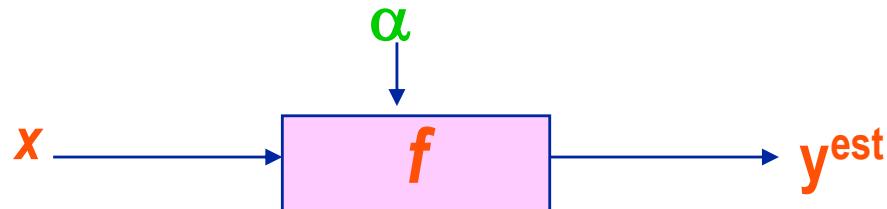
Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion 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  
School of Computer Science  
Carnegie Mellon University

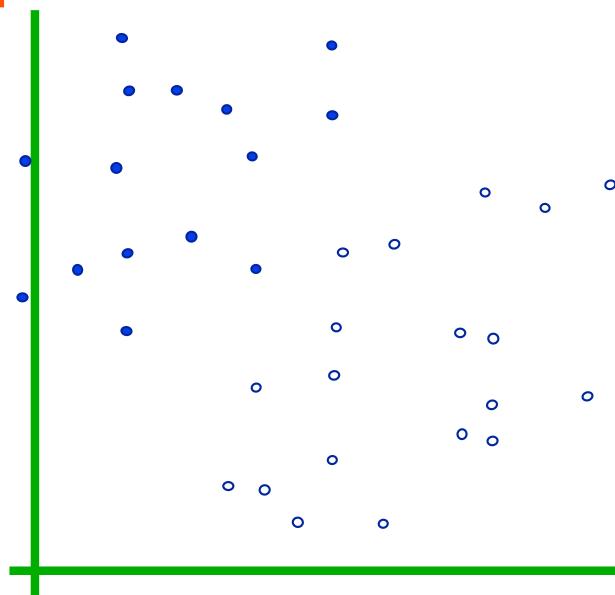
[www.cs.cmu.edu/~awm](http://www.cs.cmu.edu/~awm)  
awm@cs.cmu.edu

With minor  
modifications  
by Lyle Ungar

# Linear Classifiers



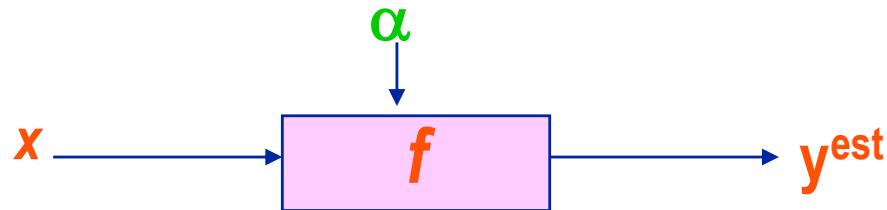
- denotes +1
- denotes -1



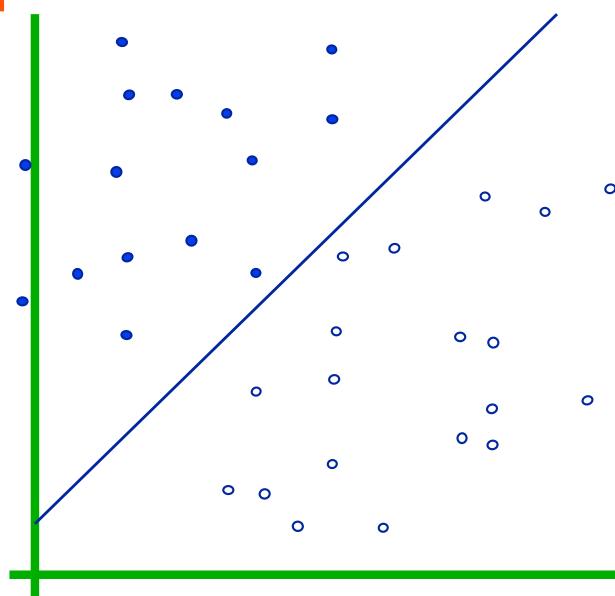
$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

How would you  
classify this  
data?

# Linear Classifiers

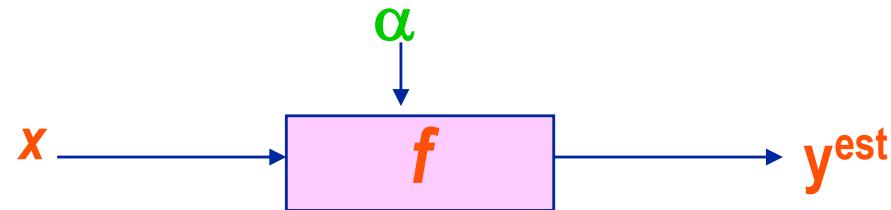


- denotes +1
- denotes -1



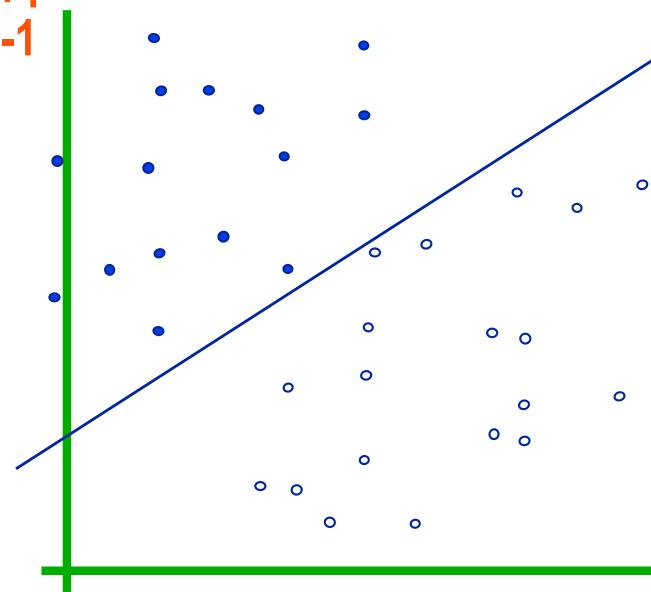
How would you  
classify this  
data?

# Linear Classifiers



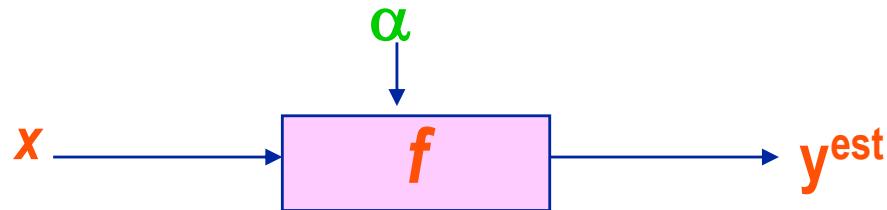
- denotes +1
- denotes -1

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

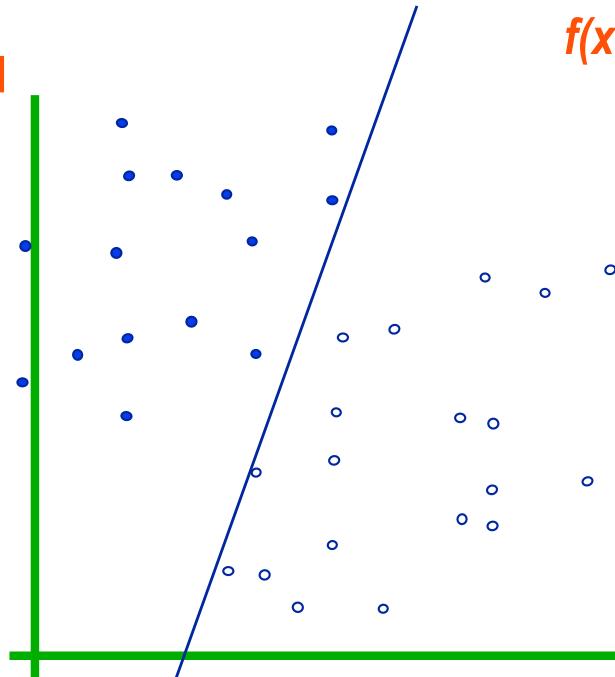


How would you  
classify this  
data?

# Linear Classifiers



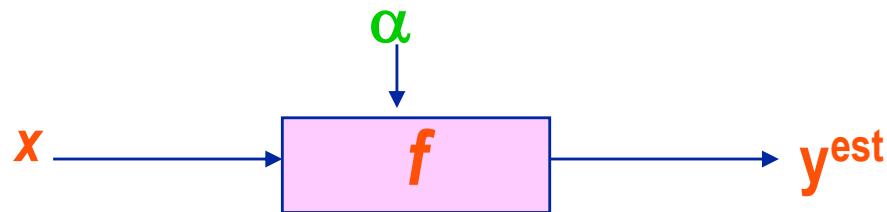
- denotes +1
- denotes -1



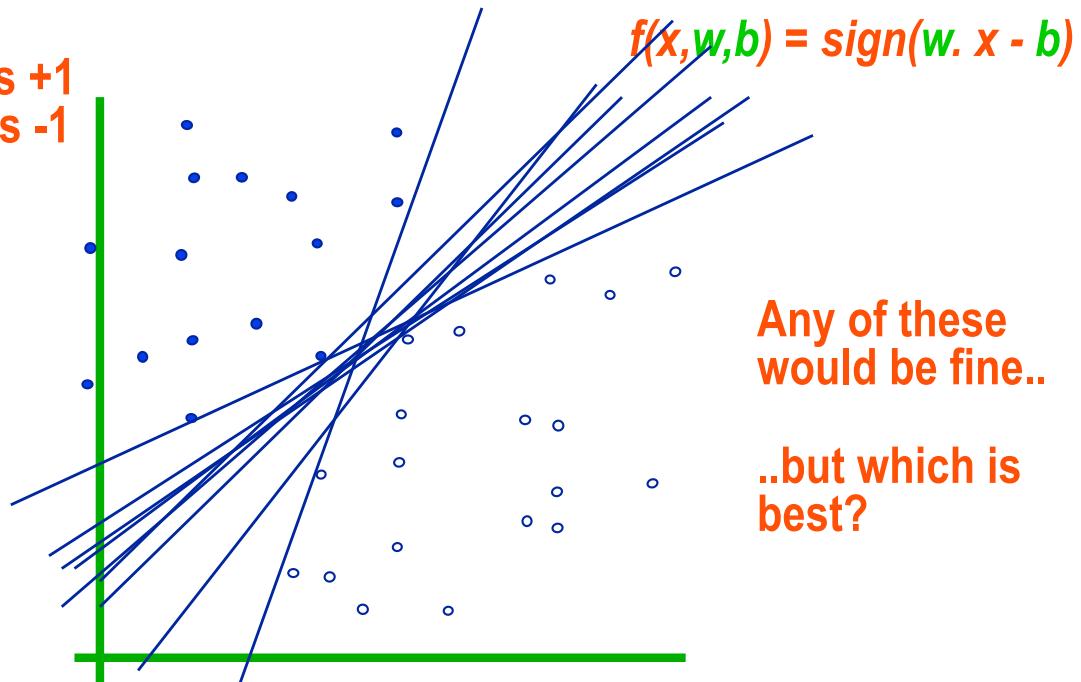
$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

How would you  
classify this  
data?

# Linear Classifiers

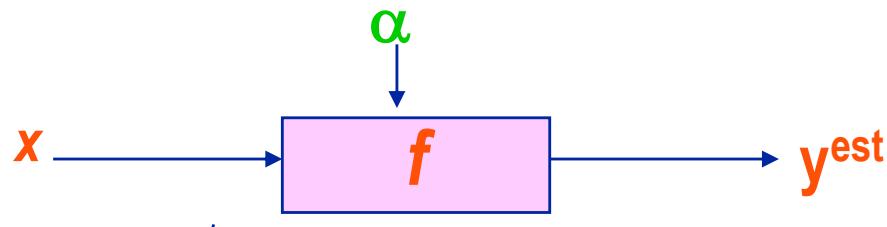
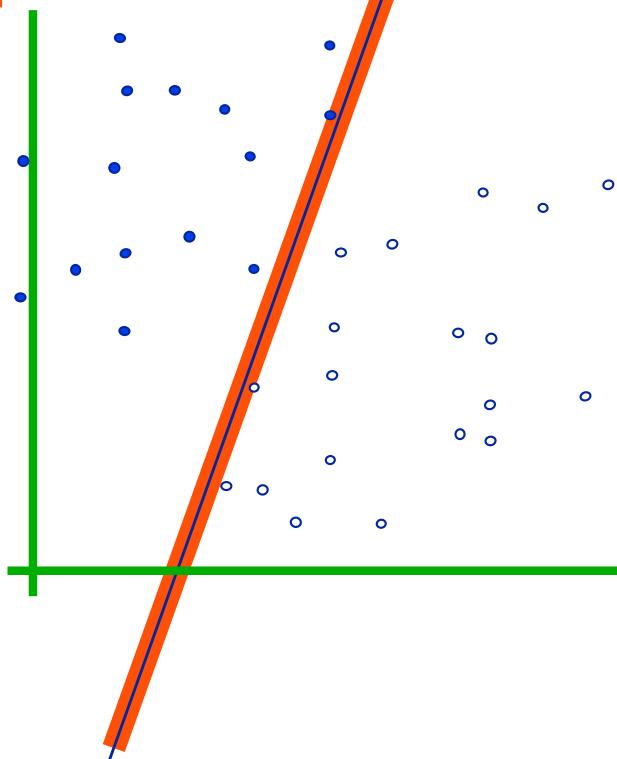


- denotes +1
- denotes -1



# Classifier Margin

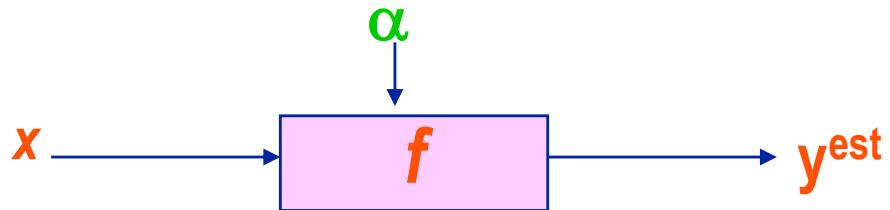
- denotes +1
- denotes -1



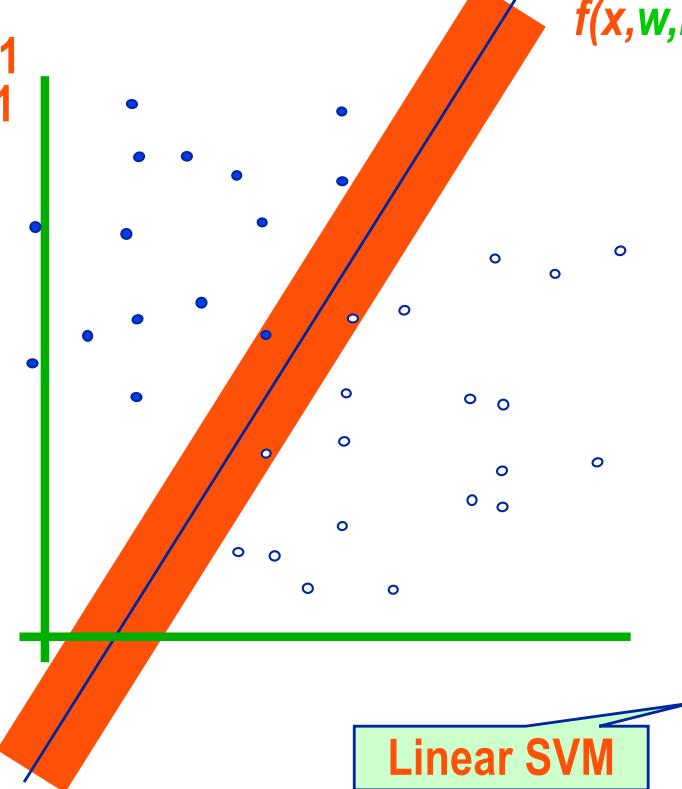
$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

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

# Maximum Margin



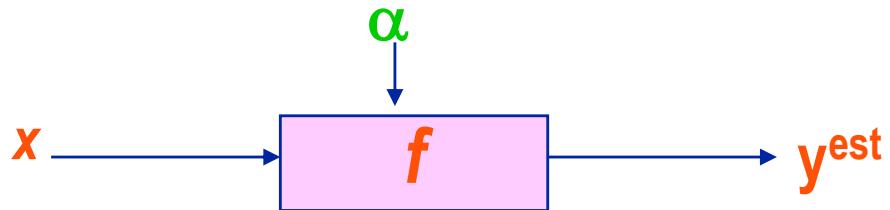
- denotes +1
- denotes -1



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

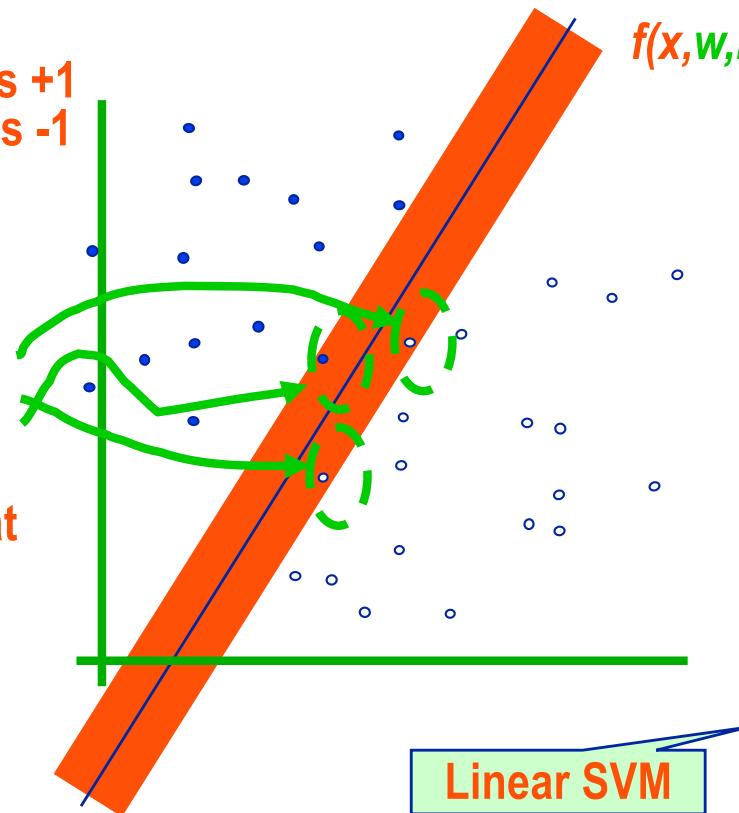
The maximum margin linear classifier is the linear classifier with the, um, maximum margin.  
This is the simplest kind of SVM (Called an LSVM)

# Maximum Margin



- denotes +1
- denotes -1

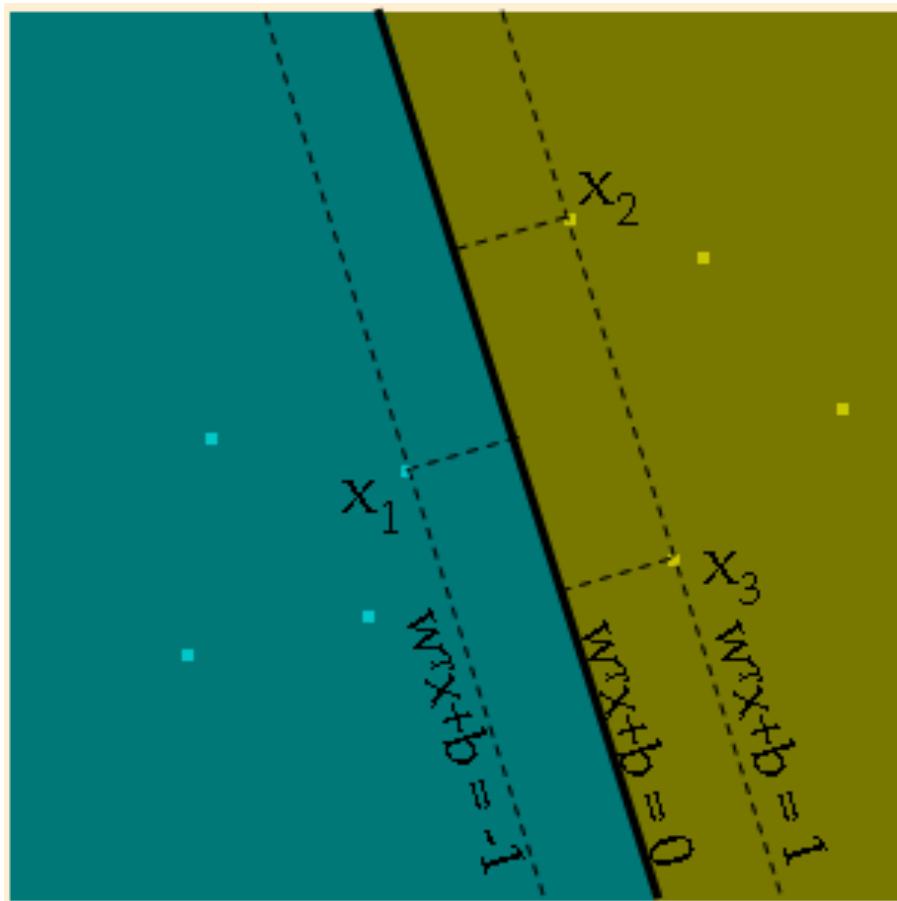
Support Vectors are those datapoints that the margin pushes up against



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

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



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

Arbitrarily normalize

$$\mathbf{w}^\top \mathbf{x} + b = \pm 1$$

$\mathbf{x}_1, \mathbf{x}_2$  and  $\mathbf{x}_3$  are support vectors

**Compute margin**

$$\mathbf{w}^\top \mathbf{x}_1 + b = -1 \quad \text{and} \quad \mathbf{w}^\top \mathbf{x}_2 + b = 1$$

**Maximize margin**

$$\mathbf{w}^\top (\mathbf{x}_2 - \mathbf{x}_1) = 2 \quad \rightarrow \quad \frac{\mathbf{w}^\top}{2\|\mathbf{w}\|_2} (\mathbf{x}_2 - \mathbf{x}_1) = \frac{1}{\|\mathbf{w}\|_2}$$

# Max margin interp. of SVM

Separable SVM primal:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^\top \mathbf{w}$$

$$\text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, n$$

# Use Lagrange Multiplier magic

Separable SVM dual:

$$\max_{\alpha \geq 0} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$$

$$\text{s.t.} \quad \sum_i \alpha_i y_i = 0$$

Most constraints are non-binding so most  $\alpha_i$  are zero.

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

The  $\mathbf{x}_i$  with nonzero  $\alpha_i$  are support vectors.

## Kernelized separable dual:

$$\max_{\alpha \geq 0} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{s.t.} \quad \sum_i \alpha_i y_i = 0$$

$$\mathbf{w}^\top \mathbf{x} + b = \sum_i \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b$$

# Max margin interp. of SVM

Separable SVM primal:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^\top \mathbf{w}$$

$$\text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, n$$

# The non-separable case

Hinge primal:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi \geq 0} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \end{aligned}$$

$$\xi_i = \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

“Slack variable” – hinge loss from the margin

# Generalize it!

Hinge primal:

$$\begin{aligned} \min_{\mathbf{w}, b, \xi \geq 0} \quad & \frac{1}{2} \|\mathbf{w}\|_p^p + C \|\xi\|_q^q \\ \text{s.t.} \quad & y_i (\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, n \end{aligned}$$

# The non-separable dual

Hinge dual:

$$\max_{\alpha \geq 0} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$$

$$\text{s.t.} \quad \sum_i \alpha_i y_i = 0, \quad \alpha_i \leq C, \quad i = 1, \dots, n$$

$\mathbf{x}_i^\top \mathbf{x}_j$  is the kernel matrix

C controls regularization

# What you should know

- ◆ Hinge loss
  - Slack variable
- ◆ Margin
- ◆ Support vector
- ◆ Primal/dual

Note: we did not cover Lagrange multipliers this year, you are not responsible for them!