

## Lecture 5: Linear Classification

Classification: Given labelled data

$$(x_i, y_i), i=1, \dots, n$$

feature vector  $\leftarrow$   $\downarrow$  label (discrete value)

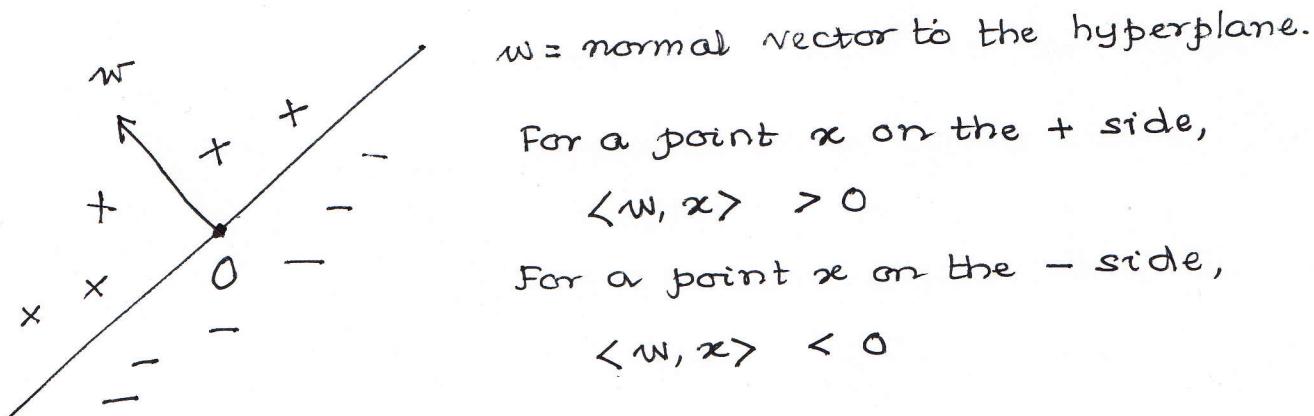
find a rule to predict  $y$  for unseen  $x$ .

For linear classification:

- assume  $y$  is  $+1$  or  $-1$
- a ~~rule~~ classification rule is a hyperplane that separates  $+1$  from  $-1$

How to represent this hyperplane?

\* For now, assume hyperplane is through origin



So :

(1) classification rule : represented by a vector  $w$

(2) Given a test example  $x$ , output its label as:

$$y = \text{sign}(\langle w, x \rangle)$$

corresponding to a hyperplane

(3) Given training data, find a  $w$  that largely separates the  $(x_i, +)$  training points on one side, and the  $(x_i, -)$  ones on another.

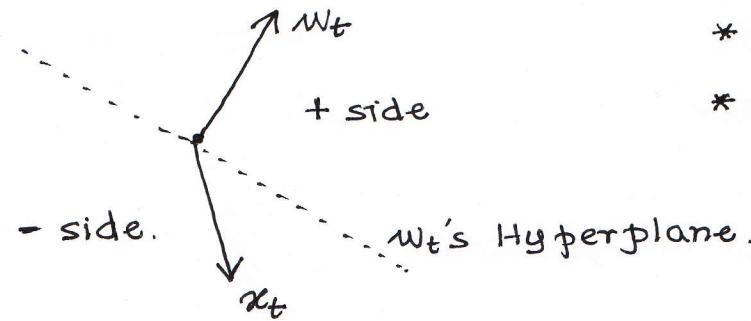
## The Perceptron Algorithm:

1. Initially,  $w_1 = 0$ .
2. For  $t = 1, 2, 3, \dots$   
 If  $y_t \langle w_t, x_t \rangle \leq 0$  then  
 $w_{t+1} = w_t + y_t x_t$   
 Else  
 $w_{t+1} = w_t$

Examples arrive in a sequence:  
 $(x_1, y_1), (x_2, y_2), \dots, (x_t, y_t)$ ,

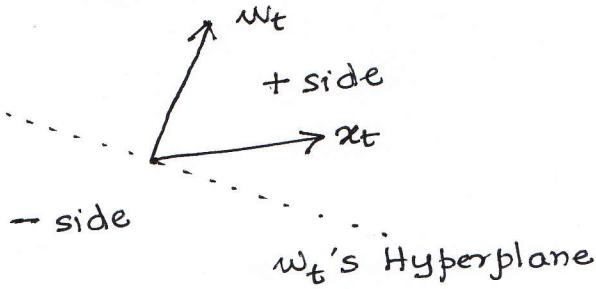
\* What does the condition " $y_t \langle w_t, x_t \rangle \leq 0$ " mean?

- Case 1:  $y_t = 1$ . Then:  $\langle w_t, x_t \rangle \leq 0$



- \*  $x_t$ 's label is +
- \* yet angle b/w  $x_t$  and  $w_t \geq 90^\circ$ ; so  $x_t$  is not on the correct side of the hyperplane corresponding to  $w_t$ .

- Case 2:  $y_t = -1$ . Then:  $\langle w_t, x_t \rangle \geq 0$



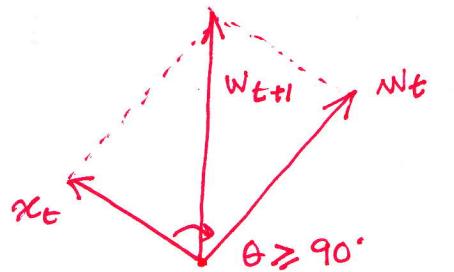
- \*  $x_t$ 's label is -
- \* yet angle b/w  $x_t$  and  $w_t$  is  $\leq 90^\circ$ , so  $x_t$  is not on the correct side of the hyperplane corresponding to  $w_t$

In both cases,  $y_t \langle w_t, x_t \rangle \leq 0$  means:

$w_t$  predicts the label of  $x_t$  incorrectly

Geometric interpretation of: " $w_{t+1} = w_t + \gamma_t x_t$ ".

Case 1:  $\gamma_t = 1$ ;  $w_{t+1} = w_t + x_t$

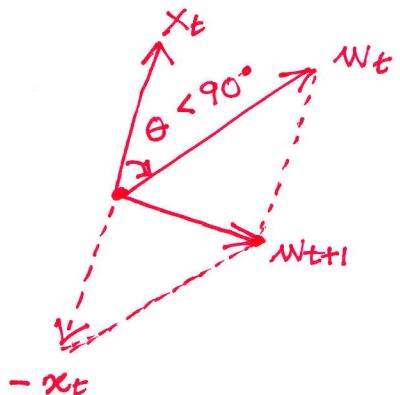


$w_{t+1}$  moves "closer to"  $x_t$

or

$x_t$  moves towards the + side of  
the decision boundary (hyperplane)

Case 2:  $\gamma_t = -1$ ;  $w_{t+1} = w_t - x_t$ .



$w_{t+1}$  moves "away from"  $x_t$

or

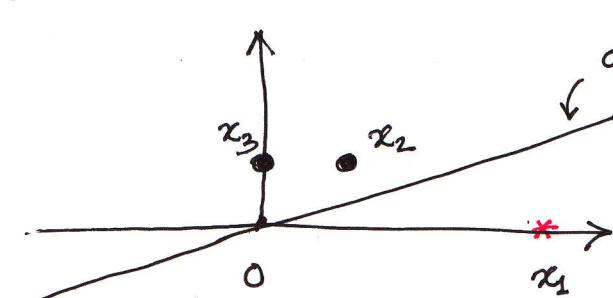
$x_t$  moves towards the - side of  
the decision boundary hyperplane.

In both cases, we are moving towards the "correct solution".

Exercise: Suppose  $w_t$  makes a mistake on  $(x_t, y_t)$ , and we update  $w_{t+1}$  as  $w_{t+1} = w_t + \gamma_t x_t$ . Is it possible for  $w_{t+1}$  to also make a mistake on  $(x_t, y_t)$ ?

### Example:

Training Data:  $((4, 0), 1), ((1, 1), -1), ((0, 1), -1), ((-2, -2), 1)$

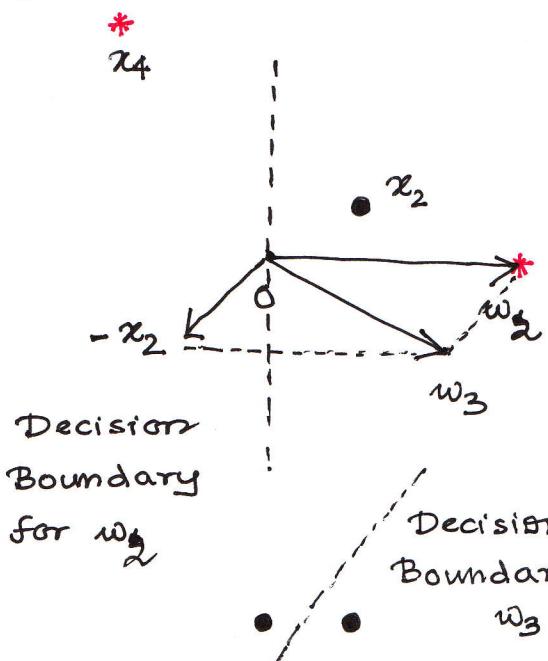


or separating  
hyperplanes  
for the input  
data points

#### Round 1:

$$* w_1 = 0$$

$$* y_t \langle w_t, x_t \rangle \text{ for } t=1 \\ = 0 \text{ as well.}$$

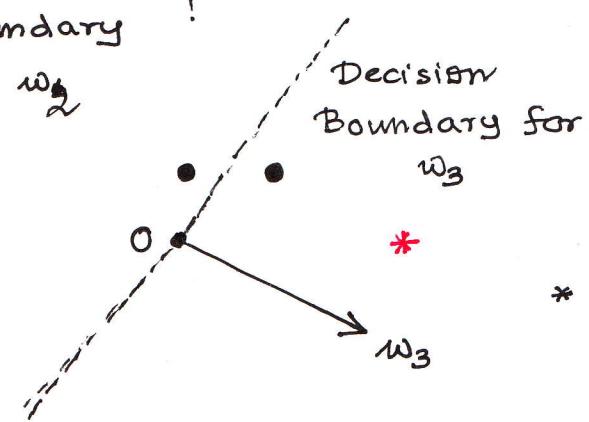


Decision  
Boundary  
for  $w_2$

#### Round 2:

$$* y_2 \langle w_2, x_2 \rangle < 0$$

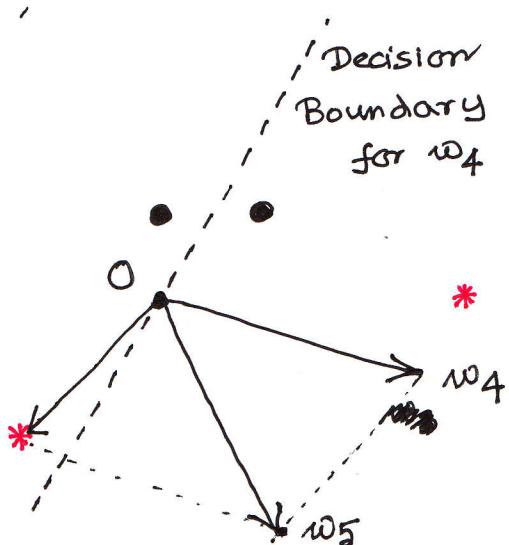
$$* \text{So } w_3 = w_2 + y_2 x_2 \\ = (4, 0) - (1, 1) = (3, -1)$$



#### Round 3:

$$* y_3 \langle w_3, x_3 \rangle > 0$$

$$\text{Correct. } w_4 = w_3 = (3, -1)$$



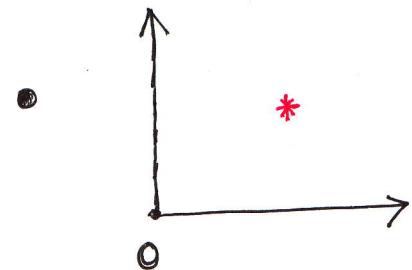
#### Round 4:

$$* y_4 \langle w_4, x_4 \rangle < 0$$

$$* \text{So } w_5 = w_4 + y_4 x_4 \\ = (3, -1) + (-2, -2) \\ = (1, -3)$$

Example 2: Training data:

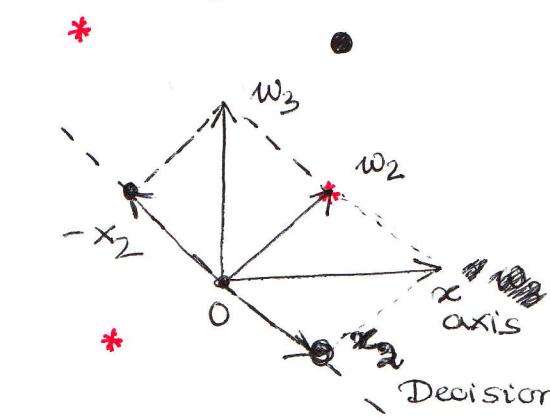
$$((1, 1), 1), ((1, -1), -1), ((-1, 1), -1), ((-1, -1), 1)$$



Initially,  $w_1 = 0$ . In round 1,

\*  $y_1 \langle w_1, x_1 \rangle \leq 0$ , so mistake.

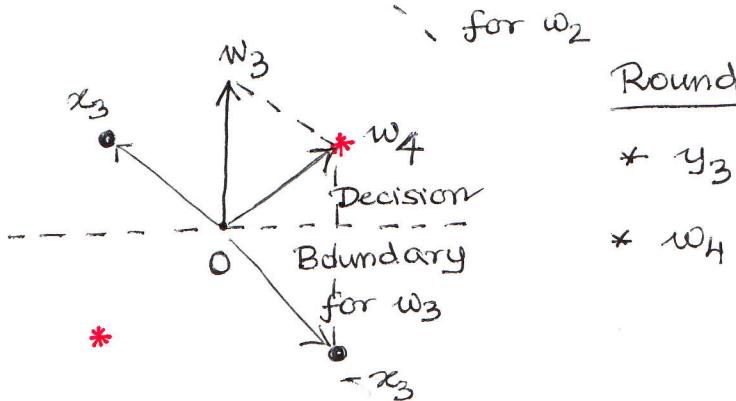
$$w_2 = w_1 + y_1 x_1 = (1, 1)$$



Round 2:

\*  $y_2 \langle w_2, x_2 \rangle \leq 0$

$$\begin{aligned} * w_3 &= w_2 + y_2 x_2 = (1, 1) - (1, -1) \\ &= (0, 2) \end{aligned}$$



Round 3:

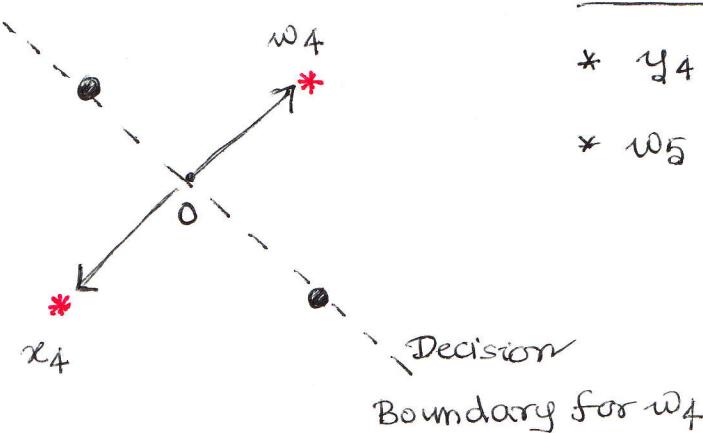
\*  $y_3 \langle w_3, x_3 \rangle \leq 0$

$$\begin{aligned} * w_4 &= w_3 + y_3 x_3 = (0, 2) - (-1, 1) \\ &= (1, 1) \end{aligned}$$

Round 4:

\*  $y_4 \langle w_4, x_4 \rangle \leq 0$ .

$$\begin{aligned} * w_5 &= w_4 + y_4 x_4 = (1, 1) + (-1, -1) \\ &= 0 \end{aligned}$$



## When does Perceptron Converge?

Linear separability: there exists a hyperplane separating the + labelled data from the -.

$\begin{matrix} + & + & + \\ - & - & - \end{matrix}$   
 Linearly Separable

$\begin{matrix} + & + \\ - & + \\ + & + \end{matrix}$   
 Not Linearly Separable

## Measure of Separability: Margin.

For a vector  $w$ , and training set  $S$ , margin of  $w$  wrt  $S$  is:

$$\gamma = \min_{(x,y) \in S} \frac{|\langle w, x \rangle|}{\|w\|}$$

Example:  $S = ((1, -1), 1), ((-1, -1), 1), ((0.01, 0), 1), ((-1, 0), -1)$

$$w = (1, 0)$$

\* What is the margin of  $w$  wrt  $S$ ?  $\|w\| = 1$ .

$$* \frac{|\langle w, x_1 \rangle|}{\|w\|} = 1$$

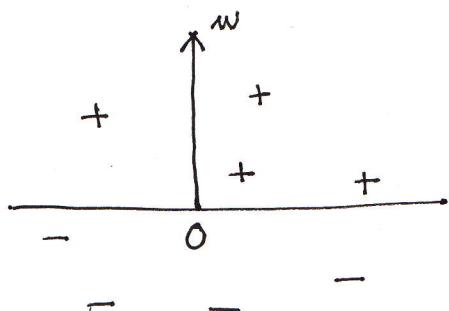
$$\text{So, margin} = 0.01$$

$$\frac{|\langle w, x_2 \rangle|}{\|w\|} = 1$$

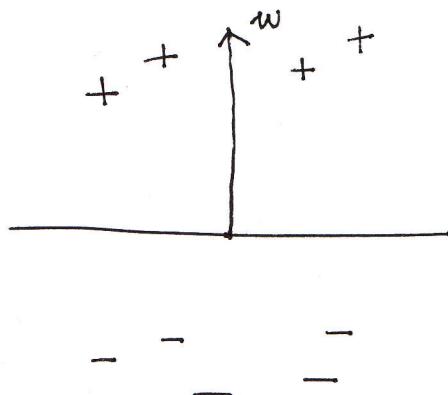
$$\frac{|\langle w, x_3 \rangle|}{\|w\|} = 0.01$$

$$\frac{|\langle w, x_4 \rangle|}{\|w\|} = 1$$

Geometrically:



Low Margin Data



High Margin Data

Theorem: If the training data is linearly separable with margin  $\gamma$ , and if  $\|x_i\| \leq 1$  for all  $i$  in the training set, then, perceptron makes  $\leq 1/\gamma^2$  mistakes.

Note:

1. Lower margin  $\Rightarrow$  more mistakes.
2. May need  $> 1$  pass over training data to get a classifier with no mistakes.

Proof: Let  $w^*$  be a linear separator with margin  $\gamma$  on the training set s.t.  $\|w^*\| = 1$ .

Fact 1: If there is a mistake at round  $t$ ,  $\langle w_{t+1}, w^* \rangle \geq \langle w_t, w^* \rangle + \gamma$

Proof: On a mistake,

$$w_{t+1} = w_t + y_t x_t. \text{ So:}$$

$$\langle w_{t+1}, w^* \rangle = \langle w_t + y_t x_t, w^* \rangle = \langle w_t, w^* \rangle + y_t \underbrace{\langle x_t, w^* \rangle}_{\geq \gamma} \geq \gamma$$

(Separable wrt  $w^* \Rightarrow y_t \langle x_t, w^* \rangle > 0$ ,

also,  $|\langle x_t, w^* \rangle| \geq \gamma$  by definition of margin)

Fact 2: If there is a mistake at round  $t$ ,

$$\|w_{t+1}\|^2 \leq \|w_t\|^2 + 1.$$

$$\begin{aligned} \|w_{t+1}\|^2 &= \|w_t + y_t x_t\|^2 = \|w_t\|^2 + y_t^2 \|x_t\|^2 + 2y_t \langle w_t, x_t \rangle \\ &\quad \underbrace{\qquad\qquad\qquad}_{\leq 1, \text{ as}} \quad \underbrace{\qquad\qquad\qquad}_{\leq 0, \text{ as there is}} \\ &\quad \|x_t\|^2 \leq 1. \quad \text{a } \underline{\text{mistake}} \text{ at} \\ &\leq \|w_t\|^2 + 1. \quad \text{round } t \end{aligned}$$

Suppose there are  $K$  mistakes.

After  $K$  mistakes,

$$\|w_{t+1}\| \leq \sqrt{K}.$$

$$\langle w_{t+1}, w^* \rangle \geq \gamma K$$

$$\text{Now: } \cos(\text{Angle b/w } w_{t+1}, w^*) = \frac{\langle w_{t+1}, w^* \rangle}{\|w^*\| \|w_{t+1}\|} \leq 1.$$

$$\text{So, } \frac{\gamma K}{\sqrt{K}} \leq 1 \Rightarrow K \leq \frac{1}{\gamma^2} \quad (\text{Proved}).$$

What if data is not linearly separable?

Ideally, we want to find a linear separator that makes the minimum number of mistakes on training data, but this is NP Hard.

But perceptron (and some other linear classification algorithms) still work when data is almost linearly separable — that is, there are a few mistakes, close to the decision boundary.

However, perceptron will never converge to a single  $w$  if the data is not linearly separable, as we make more passes over training data.

## Noted and Averaged Perception:

### Perception:

Initially,  $m = 1$ ;  $w_1 = 0$

For  $t = 1, 2, 3, \dots$

If  $y_t \langle w_m, x_t \rangle \leq 0$  then

$$w_{m+1} = w_m + y_t x_t$$

$$m = m + 1$$

Output  $w_m$

### Voted Perception:

Initially,  $m = 1$ ,  $w_1 = 0$ ,  $c_m = 1$

For  $t = 1, 2, 3, \dots$

If  $y_t \langle w_m, x_t \rangle \leq 0$  then:

$$w_{m+1} = w_m + y_t x_t$$

$$m = m + 1$$

$$c_m = 1$$

Else:

$$c_m = c_m + 1$$

Output  $(w_1, c_1), (w_2, c_2), \dots, (w_m, c_m)$

$c_m$  is a count of how long  $w_m$  survived as a classifier

### How to classify test example $x$ ?

Output: ~~sign~~ ~~count~~

$$\text{sign} \left( \sum_{i=1}^m c_i \text{sign}(\langle w_i, x \rangle) \right)$$

A problem with voted perceptron is that we have to store all the classifiers. To solve this, we use the Averaged Perceptron.

$\Rightarrow$  Averaged Perceptron uses the same algorithm as voted perceptron, but the classification rule is different.

Averaged Perceptron Classification Rule for test example  $x$ :

$$\text{Output } \operatorname{sign} \left( \sum_{i=1}^m c_i w_i, x \right)$$

Compare with voted perceptron:

$$\text{Output } \operatorname{sign} \left( \sum_{i=1}^m c_i \operatorname{sign} \langle w_i, x \rangle \right)$$

Example: If  $c_1 = c_2 = c_3 = 1$ , then:

$$\text{Averaged rule: } \operatorname{sign} (\langle w_1 + w_2 + w_3, x \rangle)$$

$$\text{Voted rule: } \underset{\text{Majority sign of}}{\cancel{\operatorname{sign} \langle w_1, x \rangle, \langle w_2, x \rangle, \langle w_3, x \rangle}}$$

### Perceptron Example:

Training data:  $((4, 0), 1), ((1, 1), -1), ((0, 1), -1), ((-2, -2), 1)$

Initially,  $w_0 = 0, m = 1, c_m = 1$

$$\text{Round 1: } y_1 \langle w_1, x_1 \rangle = 0.$$

$$\text{So: } m = 2, c_m = 1, w_2 = w_1 + y_1 x_1 = (4, 0)$$

$$\text{Round 2: } y_2 \langle w_2, x_2 \rangle \leq 0.$$

$$\text{So: } m = 3, c_m = 1, w_3 = w_2 + y_2 x_2 = (3, -1)$$

Round 3:  $y_3 \langle w_3, x_3 \rangle > 0$

so:  $c_3 = 2$ .

Round 4:  $y_4 \langle w_3, x_4 \rangle \leq 0$

so:  $c_4 = 1$ ,  $m=4$ ,  $w_4 = w_3 + y_4 x_4 = (1, -3)$

Output:  $((4, 0), 1) \quad ((3, -1), 2) \quad ((1, -3), 1)$

$$\begin{array}{ccccccc} \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ w_2 & c_2 & w_3 & c_3 & w_4 & c_4 \end{array}$$

$w_1 = (0, 0)$  (always)

$c_1 = 1$  (but its value doesn't matter)

Noted Perceptron Rule:

$$\begin{aligned} \text{sign} \left[ \text{sign}(\langle [0], x \rangle) + \text{sign}(\langle [4], x \rangle) + 2 \text{sign}(\langle [3], x \rangle) \right. \\ \left. + \text{sign}(\langle [-3], x \rangle) \right] \end{aligned}$$

$x$  = test example.

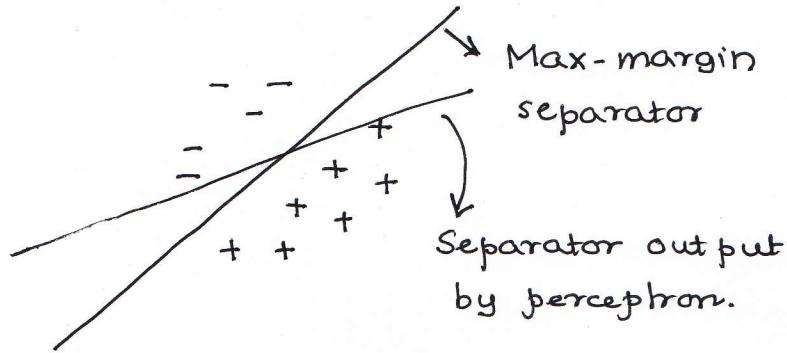
Averaged Perceptron Prediction Rule:

$$\begin{aligned} \text{sign} \left[ \langle [0] + [4] + 2 * [3] + [-3], x \rangle \right] \\ = \text{sign} \left[ \langle [11], x \rangle \right] \end{aligned}$$

Exercise: Is the voted perceptron decision boundary a hyperplane?  
What about the averaged perceptron decision boundary?

## Margins:

- Data linearly separable with a margin.
  - \* Perceptron stops once ~~a~~ a separator is found, but we may want to compute the max-margin separator.



Max-margin separators are computed by Support Vector Machines (SVMs).

## Linear Classification by a Hyperplane not thru Origin

We can transform this problem to linear classification by a hyperplane through the origin.

### Original Problem:

Training data:  $(x_i, y_i)$ ,  $i=1, \dots, n$ ,  $x_i \in \mathbb{R}^d$ ,  $y_i \in \{-1, 1\}$

Separating hyperplane:  $\alpha^T x + b \geq 0$

$(\alpha \in \mathbb{R}^d, b = \text{a scalar})$

### Transform it to:

Training data:  $(z_i, y_i)$ ,  $i=1, \dots, n$ ,  $z_i = \begin{bmatrix} 1 \\ x_i \end{bmatrix} \in \mathbb{R}^{d+1}$ ,  $y_i \in \{-1, 1\}$

Separating hyperplane:  $\alpha^T z + b \geq 0$

Separating hyperplane:  $\alpha^T z + b = 0$