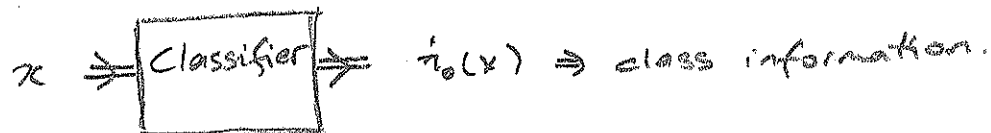


## SUPPORT VECTOR MACHINES

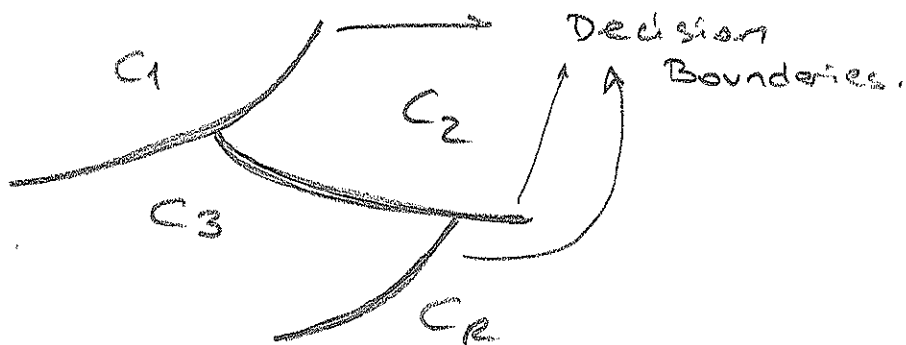
## - CLASSIFICATION

Let  $x \in \mathbb{R}^n$ , assume there are  $R$  classes  $C_1, C_2, \dots, C_m$   
 Find/build a machine/algorithm:



$i_0(x)$ : class membership function.

Assume that input space ( $\mathbb{R}^n$ ) is divided into regions.



Finding decision boundaries solves classification problem.

⇒ DISCRIMINANT FUNCTIONS

Assume that we can find  $R$  functions

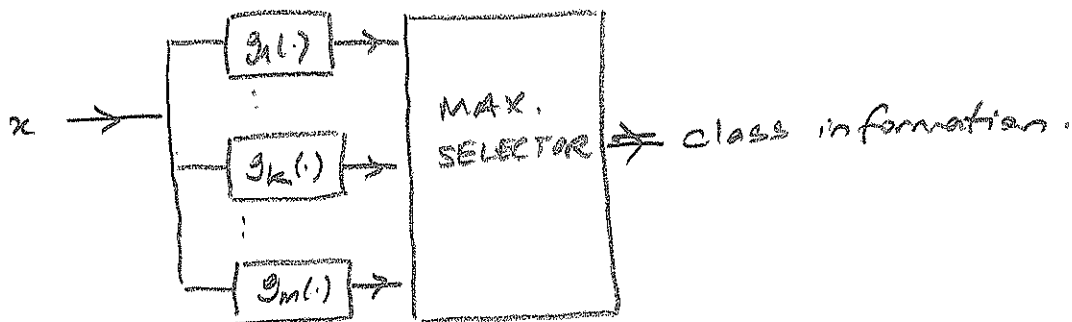
$$g_i: \mathbb{R}^n \rightarrow \mathbb{R} \quad i=1, 2, \dots$$

$$\Rightarrow x \in C_k \Leftrightarrow g_k(x) > g_j(x) \quad \begin{matrix} j=1, 2, \dots, m \\ j \neq k \end{matrix}$$

⇒ in the region of  $C_k$ , the relevant discriminant function  $g_k(\cdot)$  is maximum.

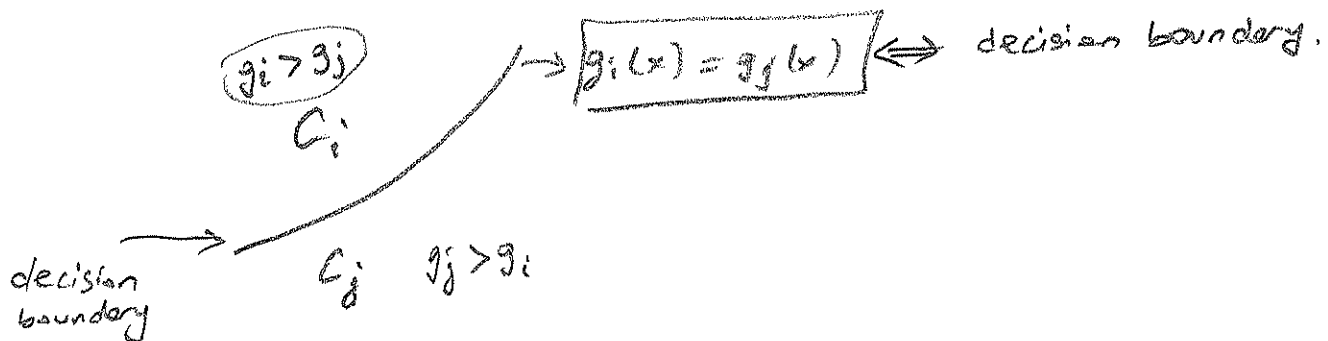
Advantages: If we can find such discriminant fns

1°) A machine can be built:



2°) If  $g_k(x)$  are LINEAR  $\Rightarrow$  a standard NN structure

3°) Decision boundaries can be computed.

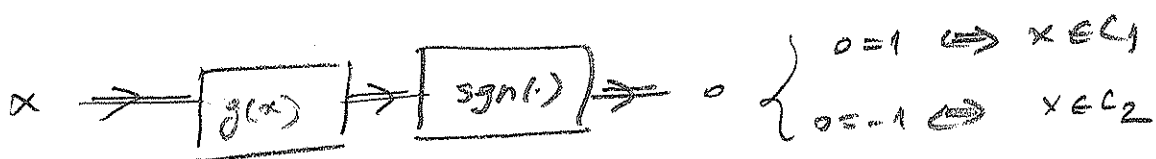


SPECIAL CASE: 2 class

$$x \in C_1 \Leftrightarrow g_1(x) > g_2(x) \Leftrightarrow g(x) = g_1(x) - g_2(x) > 0$$

$$x \in C_2 \Leftrightarrow g_2(x) > g_1(x) \Leftrightarrow g(x) = g_1(x) - g_2(x) < 0$$

$\Rightarrow$  Decision boundary  $\equiv g(x) = 0$

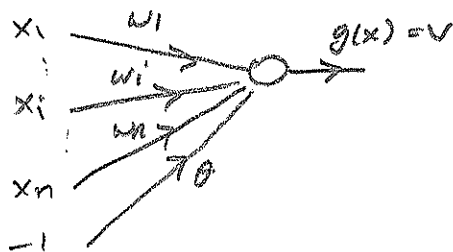


Called DICHOTOMIZER

## — LINEAR DISCRIMINANT FUNCTIONS

$$g(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n - \theta$$

⇒ related to NN

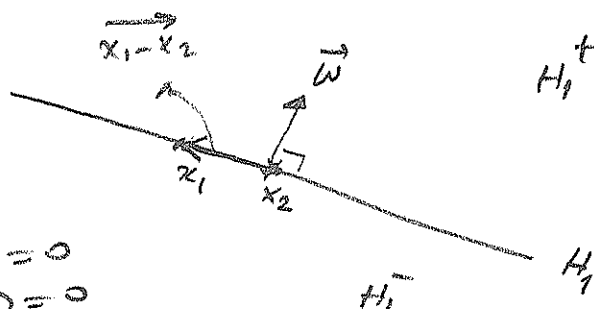


$$w = \begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$g(x) = \tilde{w}^T \tilde{x} - \theta$$

1°)  $g(x)=0$  is called a HYPERPLANE  $H_1$

2°)  $w \perp H_1$

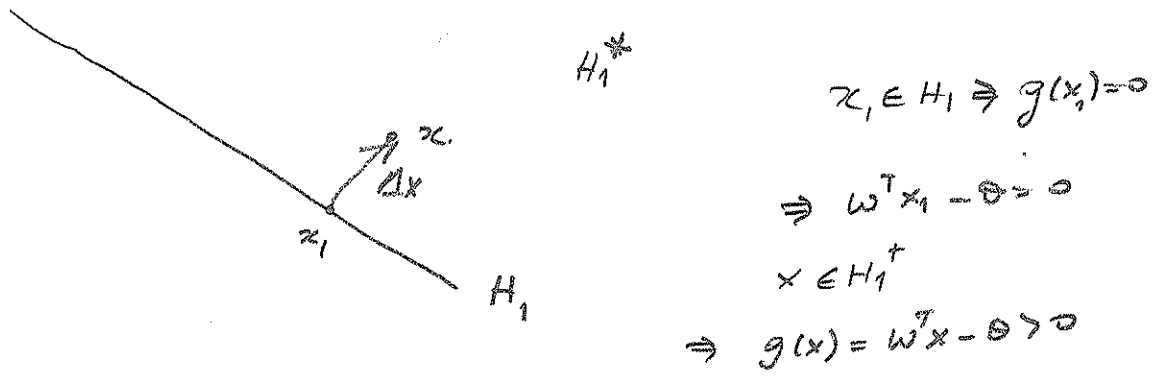


$$\text{Let } \begin{matrix} x_1 \in H_1 \\ x_2 \in H_1 \end{matrix} \Leftrightarrow \begin{matrix} g(x_1) = 0 \\ g(x_2) = 0 \end{matrix}$$

$$\Leftrightarrow \left. \begin{matrix} w^T x_1 - \theta = 0 \\ w^T x_2 - \theta = 0 \end{matrix} \right\}$$

$$w^T (x_1 - x_2) = 0$$

3°)  $w$ , as a vector points towards + side  $H_1^+$



write  $x = x_1 + \Delta x \Rightarrow g(x) = w^T x_1 - \theta + w^T \Delta x > 0$

$\Rightarrow w^T \Delta x > 0$  for any  $\Delta x$  pointing to  $H_1^+$

$\Rightarrow |\angle(w, \Delta x)| < 90^\circ$

$w^T \Delta x = \|w\| \cdot \|\Delta x\| \cos \angle(w, \Delta x) > 0$

$\Rightarrow w$  points at + side.

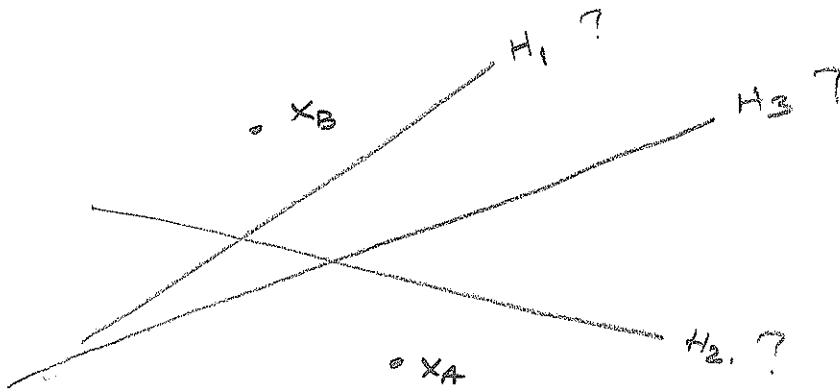
4°) Let  $x_1^*$  be on  $H_1$ ,  $x$  be a generic points on  $H_1$

$\Rightarrow \left. \begin{aligned} g(x) &= w^T x - \theta = 0 \\ g(x_1) &= w^T x_1 - \theta = 0 \end{aligned} \right\} \boxed{w^T (x - x_1) = 0}$  Line formula

$\Rightarrow g(x) = w^T x - \theta \Leftrightarrow g(x) = w^T x - w^T x_1 \Leftrightarrow \theta = w^T x_1$

(note that if  $x_1 \in H_1, x_2 \in H_1 \Rightarrow \boxed{w^T x_1 = w^T x_2}$ )

## MINIMUM DISTANCE CLASSIFIER

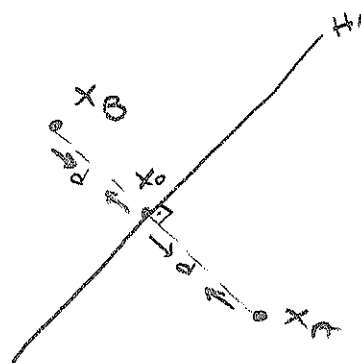


$$C_1 = \{x_B\}, \quad C_2 = \{x_A\}$$

Question: What is the BEST hyperplane separating  $C_1$  &  $C_2$ ?

Answer: The line which has EQUAL DISTANCES to both  $x_A$  and  $x_B$  is the BEST  $\Rightarrow$  May tolerate larger perturbations.

PROBLEM: Given  $x_B$  and  $x_A$ , Find the Line  $H_1$  which bisects them



$$x_0: \text{midpoint} \Leftrightarrow \text{on } H_1 \Leftrightarrow x_0 = \frac{1}{2}(x_A + x_B)$$

$$\text{normal to } H_1 \Rightarrow x_B - x_A \Rightarrow \boxed{w = (x_B - x_A)}$$

$$\text{Line formula } \boxed{w^T (x - x_0) = 0}$$

$$\boxed{(x_B - x_A)^T x - \frac{1}{2}(x_B - x_A)^T (x_B + x_A) = 0}$$

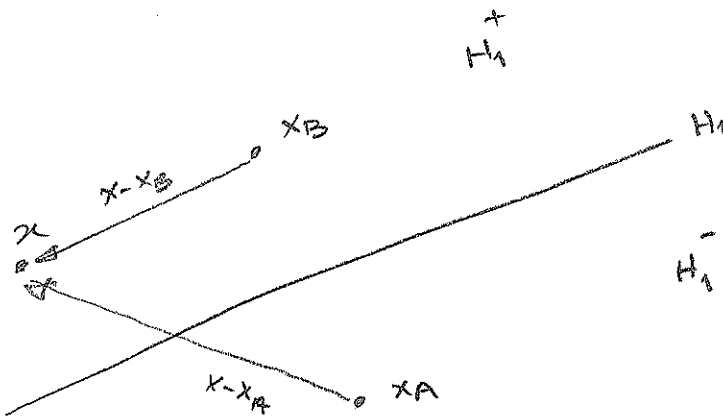
$$(x_B - x_A)^T (x_B + x_A) = x_B^T x_B - x_A^T x_A = \|x_B\|^2 - \|x_A\|^2$$

 $\Rightarrow$ 

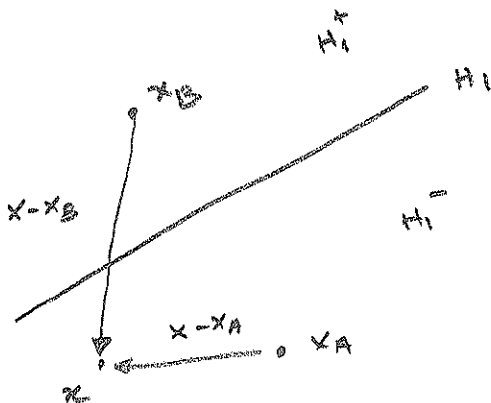
$$g(x) = w^T x - \theta$$

$$w = (x_B - x_A) \quad \theta = \|x_B\|^2 - \|x_A\|^2$$

GEOMETRICAL IDEA



$$x \in H_1^+ \Rightarrow \|x - x_B\| < \|x - x_A\|$$



$$x \in H_1^- \Rightarrow \|x - x_A\| < \|x - x_B\|$$

$$\Rightarrow x \in H_1^+ \Rightarrow \text{minimum distance to } x_B$$

$$x \in H_1^- \Rightarrow \text{minimum distance to } x_A$$

— RELATION WITH DISCRIMINANT FUNCTIONS

$$\Rightarrow g_1(x) = \|x - x_B\| \quad \text{or} \quad g_1(x) = \|x - x_B\|^2 \\ = (x - x_B)^T (x - x_B)$$

$$g_2(x) = \|x - x_A\| \quad \text{or} \quad g_2(x) = \|x - x_A\|^2 \\ = (x - x_A)^T (x - x_A)$$

$$\Rightarrow x \in C_1 \Rightarrow g_1(x) < g_2(x) \quad (\text{minimum distance})$$

$$x \in C_2 \Rightarrow g_2(x) < g_1(x)$$

$\Rightarrow$  opposite of discriminant functions

$\Rightarrow$  Actually  $-g_1(x)$  and  $-g_2(x)$  are discriminant functions.

DECISION BOUNDARY:

$$g_1(x) = g_2(x)$$

$$\Rightarrow \|x - x_B\|^2 = \|x - x_A\|^2$$

$$\Rightarrow (x - x_B)^T (x - x_B) = (x - x_A)^T (x - x_A)$$

$$\Rightarrow \|x\|^2 - 2x_B^T x + \|x_B\|^2 = \|x\|^2 - 2x_A^T x + \|x_A\|^2$$

$$\Rightarrow 2(x_B - x_A)^T x - (\|x_B\|^2 - \|x_A\|^2) = 0$$

$$\Rightarrow \boxed{(x_B - x_A)^T x - \frac{1}{2} (\|x_B\|^2 - \|x_A\|^2) = 0}$$

What if we have more than one member in classes?

$$C_1 = \{x_1, \dots, x_p\} \quad C_2 = \{x_{p+1}, \dots, x_N\}$$

SOLUTION: Choose a MEANINGFUL class representative point for each class.

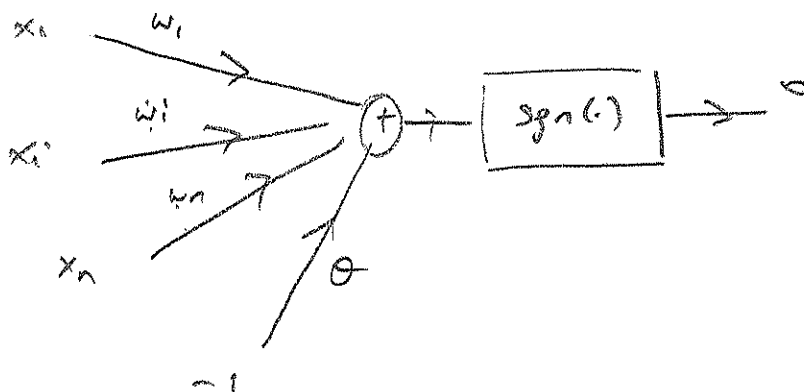
Say  $x_B$  represents  $C_1$ ,  $x_A$  represents  $C_2$

$\Rightarrow$  use min. distance classifier.

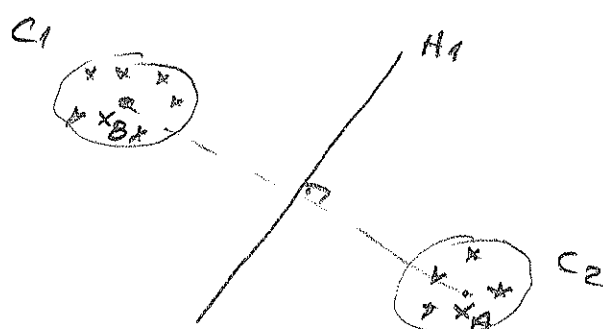
Typically,  $x_B = \text{cluster center for } C_1 = \frac{x_1 + \dots + x_p}{p}$

$$x_A = \text{cluster center for } C_2 = \frac{x_{p+1} + \dots + x_N}{N-p}$$

$$\Rightarrow w = x_B - x_A, \quad \theta = \frac{1}{2}(\|x_B\|^2 - \|x_A\|^2)$$



Works fine if  $C_1$  and  $C_2$  are WELL SEPARATED



## - MULTICATEGORY CASE

$$C_1 = \{x_1, \dots, x_5\} \rightarrow \text{class representative } \tilde{x}^1$$

$$C_2 = \{x_{5+1}, \dots, x_9\} \rightarrow \text{class representative } \tilde{x}^2$$

$$\vdots$$

$$C_i = \{x_{t+1}, \dots, x_r\} \rightarrow \text{class representative } \tilde{x}^i$$

$$\vdots$$

$$C_m = \{x_{w+1}, \dots, x_w\} \rightarrow \text{class representative } \tilde{x}^m$$

NOTE: One may choose class representatives as cluster centers.

$$x \in C_i \text{ if } \|x - x^i\|^2 < \|x - x^j\|^2 \quad j \neq i$$

$$\|x - x^i\|^2 = (x - x^i)^T (x - x^i) = \|x\|^2 - 2x^T x^i + \|x^i\|^2$$

$$\Rightarrow x \in C_i \text{ if } \|x\|^2 - 2x^T x^i + \|x^i\|^2 < \|x\|^2 - 2x^T x^j + \|x^j\|^2$$

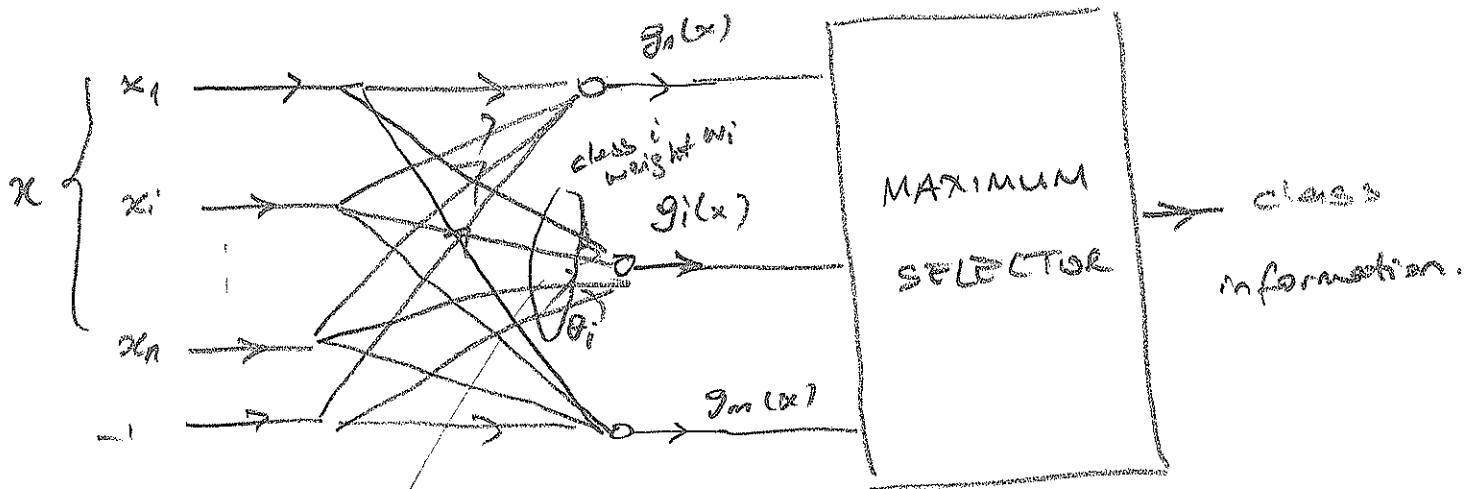
$$\Rightarrow \text{if } 2x^T x^i - \|x^i\|^2 > 2x^T x^j - \|x^j\|^2$$

$$\Rightarrow \text{Discriminant function for } C_i \Rightarrow g_i(x) = x^T x^i - \|x^i\|^2$$

$$x \in C_i \Leftrightarrow g_i(x) > g_j(x) \quad j \neq i$$

$$g_i(x) = x^T x^i - \|x^i\|^2$$

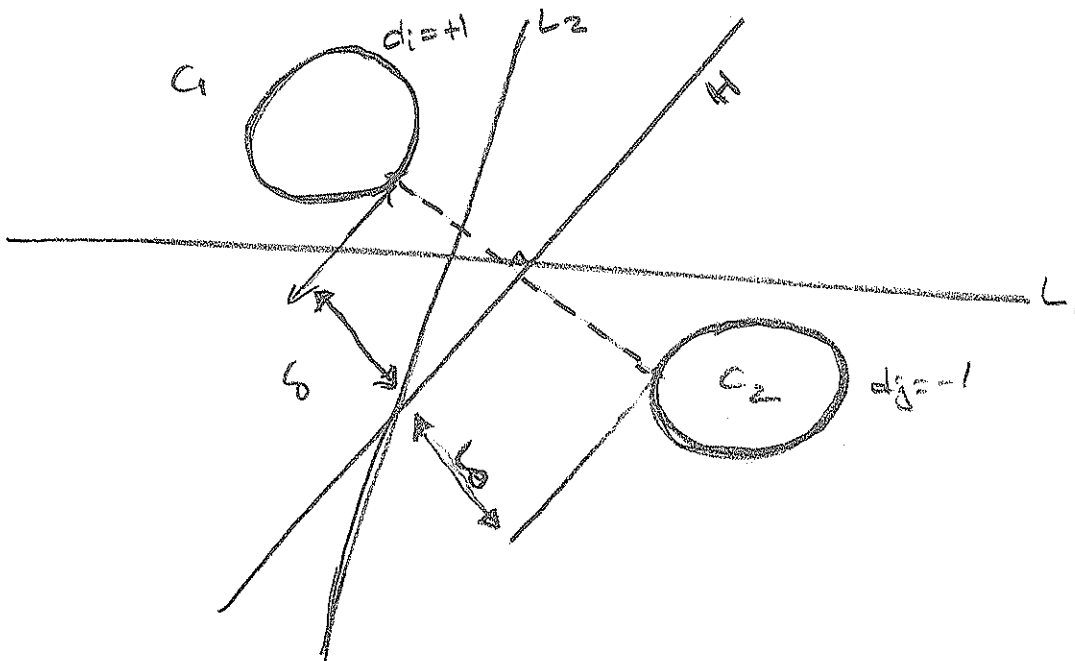
## CLASSIFIER



$$w_i = x_i \quad \theta_i = -\frac{1}{2} \|x_i\|^2$$

## SUPPORT VECTOR MACHINES (SVM)

- Extension of Minimum Distance Classifier.
- Linear case (LSVM)
- Assume that we have 2 classes
  - $C_1 = \{x_1, \dots, x_i, \dots, x_m\} \Rightarrow d_i = +1$  for  $i=1, 2, \dots, m$
  - $C_2 = \{x_{m+1}, \dots, x_j, \dots, x_N\} \Rightarrow d_j = -1$  for  $j=m+1, \dots, N$
- Assume that the classes are linearly separable.
- By using perceptron algorithm, a separating hyperplane can be found.
- Linearly separable case  $\Rightarrow$  infinitely many solutions.



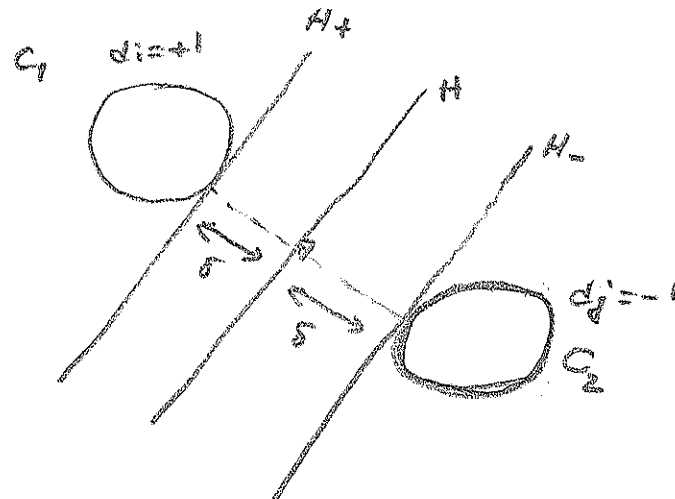
Among infinitely many solutions, which one is "the best"?

$\rightarrow$  Preferably  $H_1$ . Why?

- SVM is based on finding such a hyperplane  $H$

$$H : \boxed{g(x) = w^T x - \theta = 0}$$

- Should be on the midpoint between  $C_1$  &  $C_2$



$$H_+ \Rightarrow g_+(x) = w^T x - \theta = b$$

$$H_- \Rightarrow g_-(x) = w^T x - \theta = -b$$

$b > 0$  is a scalar.

PROPERTY 1: For simplicity, we may normalize  $b \Rightarrow 1$

Why?  $H^+ \quad w^T x - \theta = b \Leftrightarrow \frac{1}{b} w^T x - \frac{1}{b} \theta = 1$

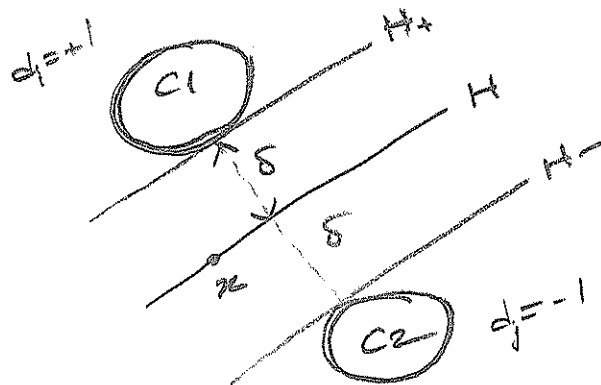
$H^- \quad w^T x - \theta = -b \Leftrightarrow \frac{1}{b} w^T x - \frac{1}{b} \theta = -1$

Define  $\hat{w} = \frac{1}{b} w \quad \hat{\theta} = \frac{1}{b} \theta$

$$\Rightarrow \boxed{\begin{aligned} H : g(x) &\Rightarrow \hat{w}^T x - \hat{\theta} = 0 \\ H_+ : g_+(x) &\Rightarrow \hat{w}^T x - \hat{\theta} = +1 \\ H_- : g_-(x) &= \hat{w}^T x - \hat{\theta} = -1 \end{aligned}}$$

- What is the distance  $\delta$  to  $C_1$  &  $C_2$ ?

- Let  $x$  be on  $H \Rightarrow \boxed{w^T x - \theta = 0}$



$\vec{w}$ : normal to  $H$

$\Rightarrow y = x + \delta w$  should be on  $H_+$

$\Rightarrow w^T y - \theta = 1 \Rightarrow \boxed{w^T x + \delta w^T w - \theta = 1}$

But  $x \in H \Rightarrow \boxed{w^T x - \theta = 0} \Rightarrow \delta w^T w = 1 \Rightarrow \boxed{\delta = \frac{1}{\|w\|^2}}$

FACT: Distance  $\delta$  is inversely proportional to  $\|w\|$ .

AIM: Maximize the distance  $\delta \Rightarrow$  minimize  $\|w\|$

PROBLEM 1:

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|w\|^2 \\ & \text{subject to} && w^T x_i - \theta \geq 1 \quad \text{if } x_i \in C_1 \quad (d_i = +1) \\ & && w^T x_j - \theta \leq -1 \quad \text{if } x_j \in C_2 \quad (d_j = -1) \end{aligned}$$

$\Rightarrow$  QUADRATIC OPTIMIZATION PROBLEM with linear constraints.

$\Rightarrow$  CONVEX  $\Rightarrow$  HAS a GLOBAL MIN. SOLUTION (UNIQUE!)

- SECOND FORM :

P2 :

$$\begin{array}{ll} \text{minimize} & \frac{1}{2} \|w\|^2 \\ \text{s.t.} & d_i (w^T x_i - \theta) \geq 1 \quad i=1, 2, \dots, N \end{array}$$

- Various solvers solve this kind of problems.

- One analytic method: LAGRANGE MULTIPLIERS

Solution of P2 is an EXTREMUM POINT of

$$P3 \quad L(w, \theta, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i [d_i (w^T x_i - \theta) - 1]$$

Here,  $\alpha_i \geq 0$  are called LAGRANGE MULTIPLIERS

NOTE: Actually, we should use min. wrt  $w, \theta$  and max wrt  $\alpha_i$ !

$$P4 : \quad \min_{w, \theta} \max_{\alpha} L(w, \theta, \alpha)$$

Why? Otherwise choose  $\alpha_i = \infty \Rightarrow$  get  $-\infty$ , which is minimum.

SOLUTION: Karush-Kuhn-Tucker (KKT) conditions.

$$\frac{\partial L}{\partial w} = 0 \Rightarrow \boxed{w - \sum_{i=1}^N \alpha_i d_i x_i = 0} \quad (1)$$

$$\frac{\partial L}{\partial \theta} = 0 \Rightarrow \boxed{\sum_{i=1}^N \alpha_i d_i = 0} \quad (2)$$

$$\frac{\partial L}{\partial \alpha_i} = 0 \Rightarrow \boxed{d_i (w^T x_i - \theta) - 1 = 0} \quad (3)$$

+ Remember:  $\boxed{\alpha_i \geq 0} \quad (4)$

$$(1) \Rightarrow \boxed{w = \sum_{i=1}^N \alpha_i d_i x_i} \quad (**)$$

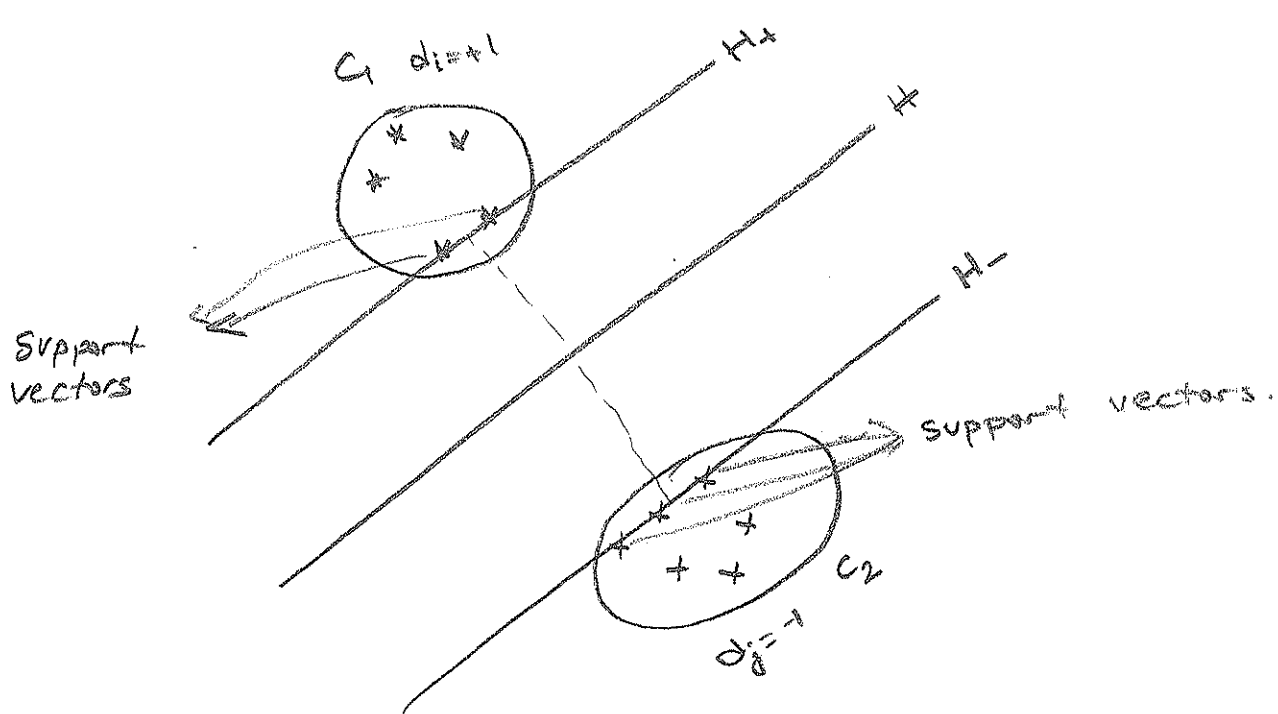
FACT:  $w$  is a linear combination of some training patterns  $x_i$ ? Which ones?

$$(3) \Rightarrow d_i (w^T x_i - \theta) - 1 = 0$$

$$d_i = +1 \Rightarrow w^T x_i - \theta = 1 \Rightarrow x_i \text{ is ON } H_+$$

$$d_i = -1 \Rightarrow w^T x_i - \theta = -1 \Rightarrow x_i \text{ is ON } H_-$$

$\Rightarrow$  The only relevant  $x_i$ 's on ~~(\*)~~ are the patterns which are on  $H_+$  and  $H_-$ . They are called SUPPORT VECTORS.



CONCLUSION  $\alpha_i > 0$  for support vectors, and  $\alpha_i = 0$  for all other patterns.

CONCLUSION: Optimal weight  $w$  is determined by only SUPPORT VECTORS. The threshold term  $\theta$  can be found easily:

if  $x_i$  is a support vector

$$\Rightarrow d_i [w^T x_i - \theta] - 1 = 0 \Rightarrow \boxed{w^T x_i - \theta = 1/d_i}$$

But  $d_i = \pm 1 \Rightarrow 1/d_i = d_i \Rightarrow w^T x_i - \theta = d_i$

$$\Rightarrow \boxed{\theta = w^T x_i - d_i}$$

Sometimes an average over the support vectors is used to find  $\theta$  for ROBUSTNESS, i.e. let  $N_s$  denote the # of support vectors

$$\boxed{\theta_{ave} = \frac{1}{N_s} \sum_{i=1}^{N_s} (w^T x_i - d_i) \quad x_i \text{ is SV}}$$

## - DUAL FORM OF OPTIMIZATION

- Since  $w = \sum_{i=1}^N \alpha_i d_i x_i$ ,  $\sum_{i=1}^N \alpha_i d_i = 0$

$\Rightarrow$  use it on  $L(w, \theta, \alpha)$

$$\begin{aligned} \Rightarrow L(w, \theta, \alpha) &= \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i [d_i (w^T x_i - \theta) - 1] \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j x_i^T x_j - \sum_{i=1}^N \alpha_i \alpha_j d_i d_j x_i^T x_j + \theta \sum_{i=1}^N \alpha_i d_i + \sum_{i=1}^N \alpha_i \end{aligned}$$

$$\Rightarrow \boxed{L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \alpha_i \alpha_j d_i d_j x_i^T x_j}$$

DUAL PROBLEM:

$$\begin{array}{l} \text{D.P.} \\ \max_{\alpha} L(\alpha) \\ \text{s.t. } \alpha_i \geq 0 \end{array}$$

$\Rightarrow$  much simpler in form  $\Rightarrow$  easy to solve.

$\Rightarrow$  various software packages.

NOTE: In the solution of D.P., some  $\alpha_i = 0$ , some  $\alpha_j > 0$ .

The ones  $\alpha_j > 0$  corresponds to SUPPORT VECTORS.

$$\Rightarrow \boxed{w = \sum_{j=1}^{N_s} \alpha_j d_j x_j}$$

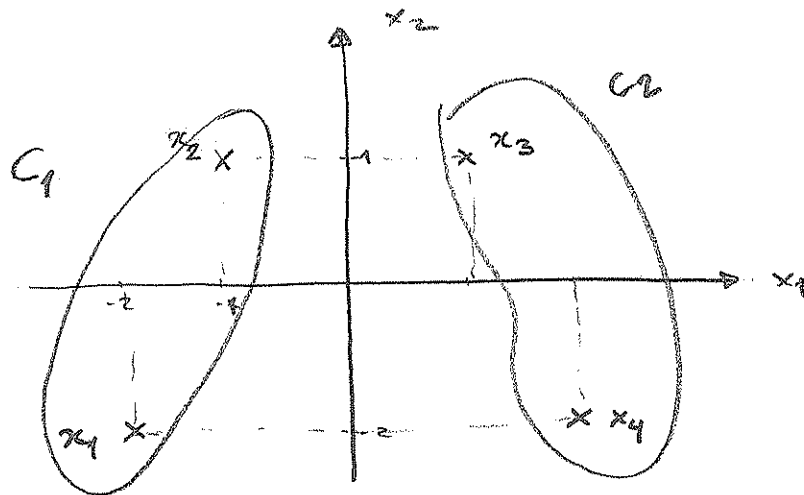
$x_j$  support vector  
 $N_s$ : # of support vectors.

$\theta$  can be found as before.

- AN EXAMPLE

$$C_1 = \left\{ x_1 = \begin{pmatrix} -2 \\ -2 \end{pmatrix}, x_2 = \begin{pmatrix} -1 \\ +1 \end{pmatrix} \right\} \quad d_1 = d_2 = +1$$

$$C_2 = \left\{ x_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, x_4 = \begin{pmatrix} 2 \\ -2 \end{pmatrix} \right\} \quad d_3 = d_4 = -1$$



ANALYTIC SOLUTION FOR SVM

$$\text{Dual Prob} \Rightarrow L(\alpha) = \sum_{i=1}^4 \alpha_i - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \alpha_i \alpha_j d_i d_j x_i^T x_j$$

$$\Rightarrow L(\alpha) = \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 - 4\alpha_1^2 - \alpha_2^2 - \alpha_3^2 - 4\alpha_4^2 - 4\alpha_1\alpha_3 - 4\alpha_2\alpha_4$$

PROBLEM:

$$\max L(\alpha)$$

$$\text{s.t. } \alpha_i \geq 0$$

- Should use quadratic optimization packages
- Or use K-T-H conditions use Linear Programming
- For this special case:

$$L(\alpha) = \underbrace{[\alpha_1 + \alpha_3 - 4\alpha_1^2 - \alpha_3^2 - 4\alpha_1\alpha_3]}_{L_1} + \underbrace{[\alpha_2 + \alpha_4 - \alpha_2^2 - 4\alpha_4^2 - 4\alpha_2\alpha_4]}_{L_2}$$

$L_1(\alpha_1, \alpha_3) \rightarrow \max$ , separately

$L_2(\alpha_2, \alpha_4) \rightarrow \max$ , separately

$$\begin{aligned} \frac{\partial L_1}{\partial \alpha_1} = 1 - 8\alpha_1 - 4\alpha_3 = 0 &\Rightarrow 2\alpha_1 + \alpha_3 = 1/4 \\ \frac{\partial L_1}{\partial \alpha_3} = 1 - 2\alpha_3 - 4\alpha_1 = 0 &\Rightarrow 2\alpha_1 + \alpha_3 = 1/2 \end{aligned} \quad \left. \vphantom{\begin{aligned} \frac{\partial L_1}{\partial \alpha_1} = 1 - 8\alpha_1 - 4\alpha_3 = 0 \\ \frac{\partial L_1}{\partial \alpha_3} = 1 - 2\alpha_3 - 4\alpha_1 = 0 \end{aligned}} \right\} \text{conflict!}$$

choose first  $\Rightarrow \alpha_3 = \frac{1}{4} - 2\alpha_1$

$$\begin{aligned} \Rightarrow \text{put in } L_1 &\Rightarrow L_1 = \alpha_1 + \left(\frac{1}{4} - 2\alpha_1\right) - 4\alpha_1^2 - \left(\frac{1}{4} - 2\alpha_1\right)^2 - 4\alpha_1\left(\frac{1}{4} - 2\alpha_1\right) \\ &= \frac{3}{16} - \alpha_1 \end{aligned}$$

For  $L_1$  max  $\Rightarrow \alpha_1 = 0 \Rightarrow \alpha_3 = \frac{1}{4} \Rightarrow L_1 = \frac{3}{16}$  ←

choose second  $\Rightarrow \alpha_3 = \frac{1}{2} - 2\alpha_1$

$$\begin{aligned} \Rightarrow \text{put in } L_1 &L_1 = \alpha_1 + \left(\frac{1}{2} - 2\alpha_1\right) - 4\alpha_1^2 - \left(\frac{1}{2} - 2\alpha_1\right)^2 - 4\alpha_1\left(\frac{1}{2} - 2\alpha_1\right) \\ &= \frac{1}{4} - \alpha_1 \end{aligned}$$

For  $L_1$  max  $\Rightarrow \alpha_1 = 0 \Rightarrow \alpha_3 = 1/2 \Rightarrow L_1 = \frac{1}{4}$  ←

Second one is bigger  $\Rightarrow$

$$\boxed{\begin{matrix} \alpha_1 = 0 \\ \alpha_3 = 1/2 \end{matrix}}$$

For the  $L_2 \rightarrow$  by symmetry  $\alpha_1 \leftrightarrow \alpha_4$   
 $\alpha_2 \leftrightarrow \alpha_3$

$$\Rightarrow \begin{cases} \alpha_4 = 0 \\ \alpha_2 = 1/2 \end{cases}$$

RESULT 1:

$x_2$  and  $x_3$  are SUPPORT VECTORS!

RESULT 2:

$$w = \alpha_2 d_2 x_2 + \alpha_3 d_3 x_3$$

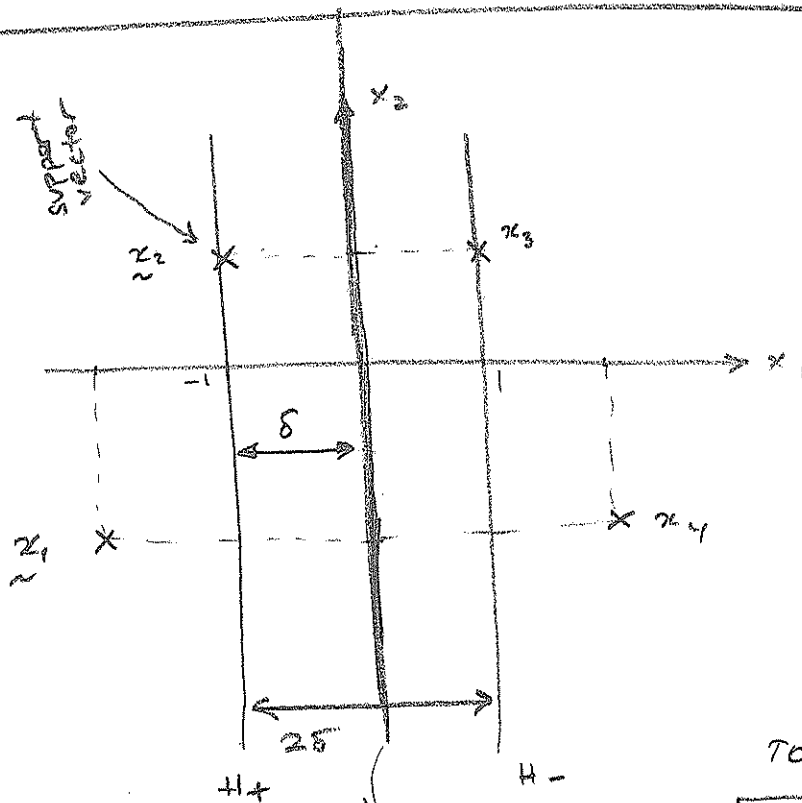
$$= \frac{1}{2} (x_2 - x_3) = \frac{1}{2} \left( \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) = \frac{1}{2} \begin{bmatrix} -2 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

$$\theta = w^T x_2 - d_2 = w^T x_3 - d_3 = 0$$

$$\Rightarrow g(x) = \tilde{w}^T \tilde{x} - \theta = -x_1 = 0 \Rightarrow \text{VERTICAL LINE.}$$

$$H_+ : w^T x - \theta = 1 \Rightarrow -x_1 = 1 \Rightarrow \boxed{x_1 = -1}$$

$$H_- : w^T x - \theta = -1 \Rightarrow -x_1 = -1 \Rightarrow \boxed{x_1 = 1}$$



TOTAL DISTANCE

$$\delta_T = 2\delta = \frac{2}{\|w\|} = 2$$

$$\delta = \frac{1}{\|w\|} = \frac{1}{1} = 1$$