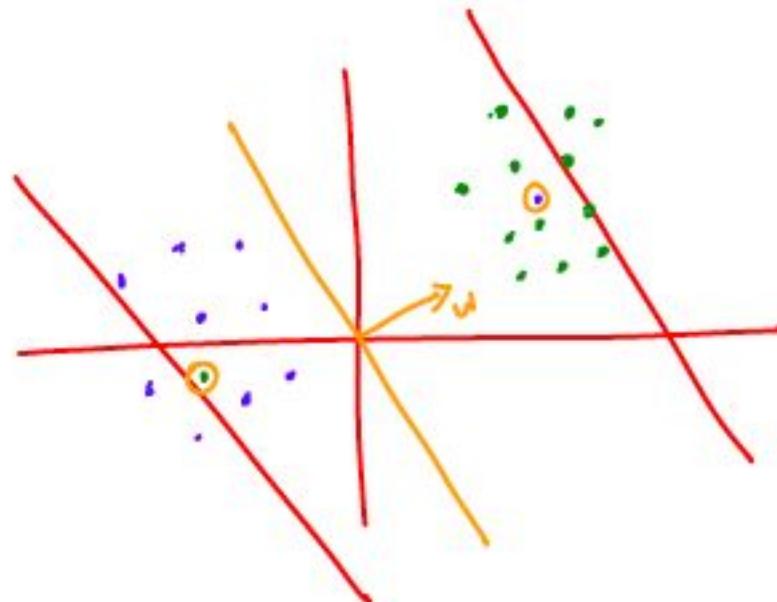
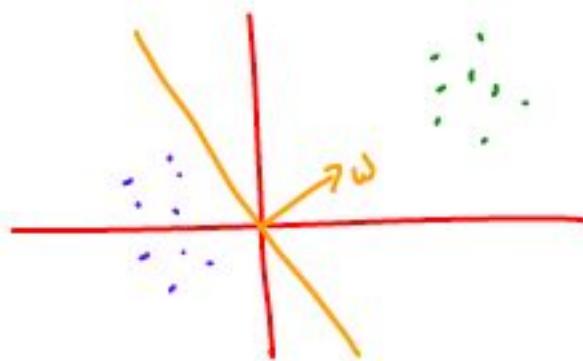


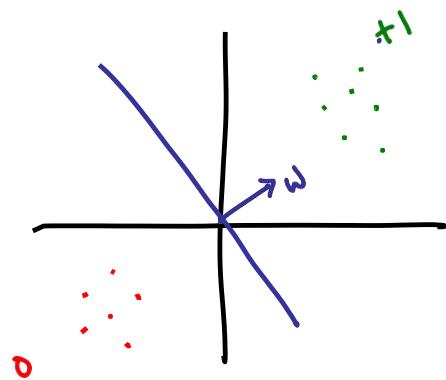
LINEAR  
SEPARABILITY  
ASSUMPTION



Dataset is  
not allowed in  
our model



Allowed.



$$P(y=1|x) = \begin{cases} 1 & \text{if } \mathbf{w}^\top \mathbf{x} \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

Linear separability assumption.

Lin. Sep. assumption :

$$\exists \mathbf{w} \in \mathbb{R}^d \text{ s.t. } \underbrace{\text{Sign}(\mathbf{w}^\top \mathbf{x}_i)}_{\downarrow} = y_i \quad \forall i \in [n] \quad \{1, \dots, n\}.$$

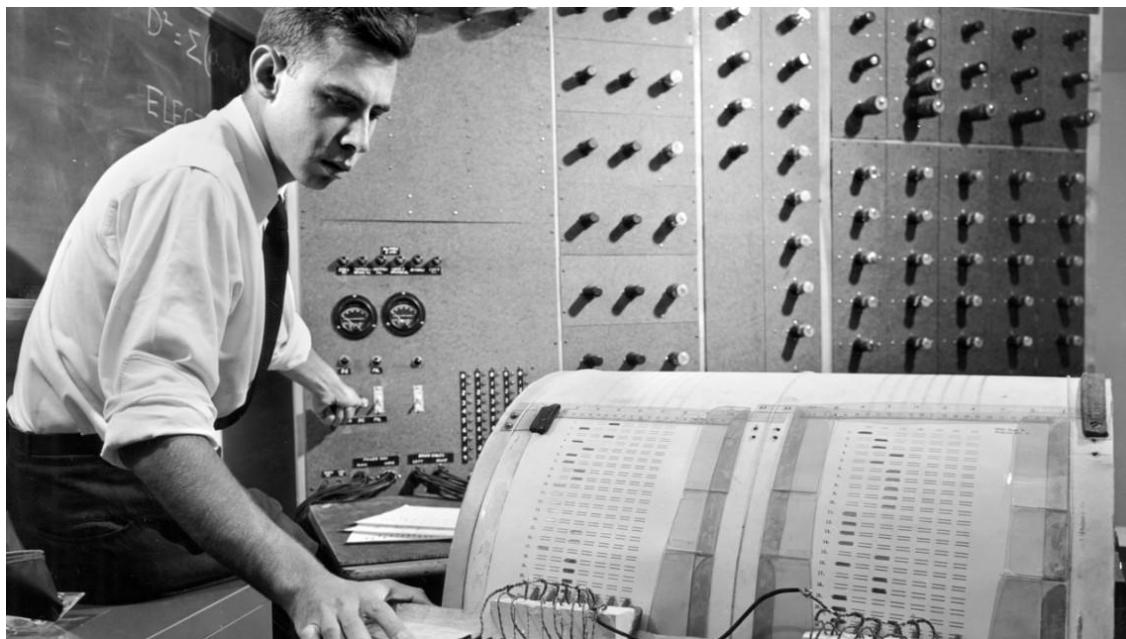
$$\text{Sign}(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

$$\begin{array}{c}
 \text{DATA} \quad \left\{ (x_1, y_1), \dots, (x_n, y_n) \right\} \\
 \min_{\underline{w \in \mathbb{R}^d}} \sum_{i=1}^n \mathbb{1}(\text{sign}(w^T x_i) \neq y_i) \rightarrow \underline{\text{NP-HARD}}
 \end{array}$$

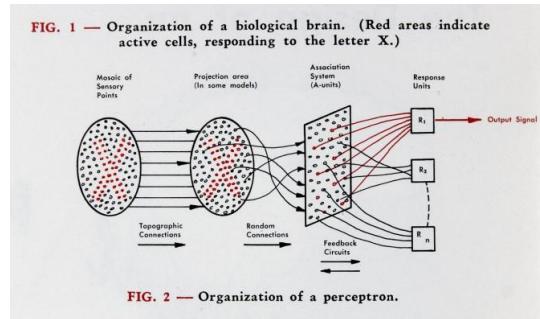
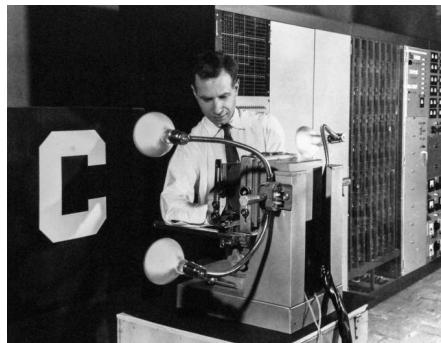
We know in general the problem of finding best  $w$  is NP-hard.

But is it still hard under linear separability assumption?

# PERCEPTRON



Frank Rosenblatt '50, Ph.D. '56, works on the "perceptron" – what he described as the first machine "capable of having an original idea."



## PERCEPTRON - ALGORITHM

Input:  $\{(x_1, y_1), \dots, (x_n, y_n)\}$   $x_i \in \mathbb{R}^d$   
 $y_i \in \{+1, -1\}$

$w^0 = [0 \ 0 \ \dots \ 0]$   $o \in \mathbb{R}^d$

until convergence

→ Pick  $(x_i, y_i)$  from dataset

IF  $(\text{sign}(w^t \cdot x_i) = y_i)$

do nothing

else

$$w^{t+1} = w^t + x_i y_i \in \mathbb{R}^d$$

UPDATE  
RULE.

end.

end

### UPDATE RULE

$$w^{t+1} = w^t + x_i y_i$$



Two types of mistake

Type 1

$$\begin{aligned} \text{Pred} &\rightarrow 1 \\ \text{act} &\rightarrow -1 \end{aligned}$$

Type 2

$$\begin{aligned} \text{Pred} &\rightarrow -1 \\ \text{act} &\rightarrow 1 \end{aligned}$$

Type -1

$$\begin{aligned} \text{Pred} &\rightarrow 1 \\ \text{act} &\rightarrow -1 \end{aligned}$$

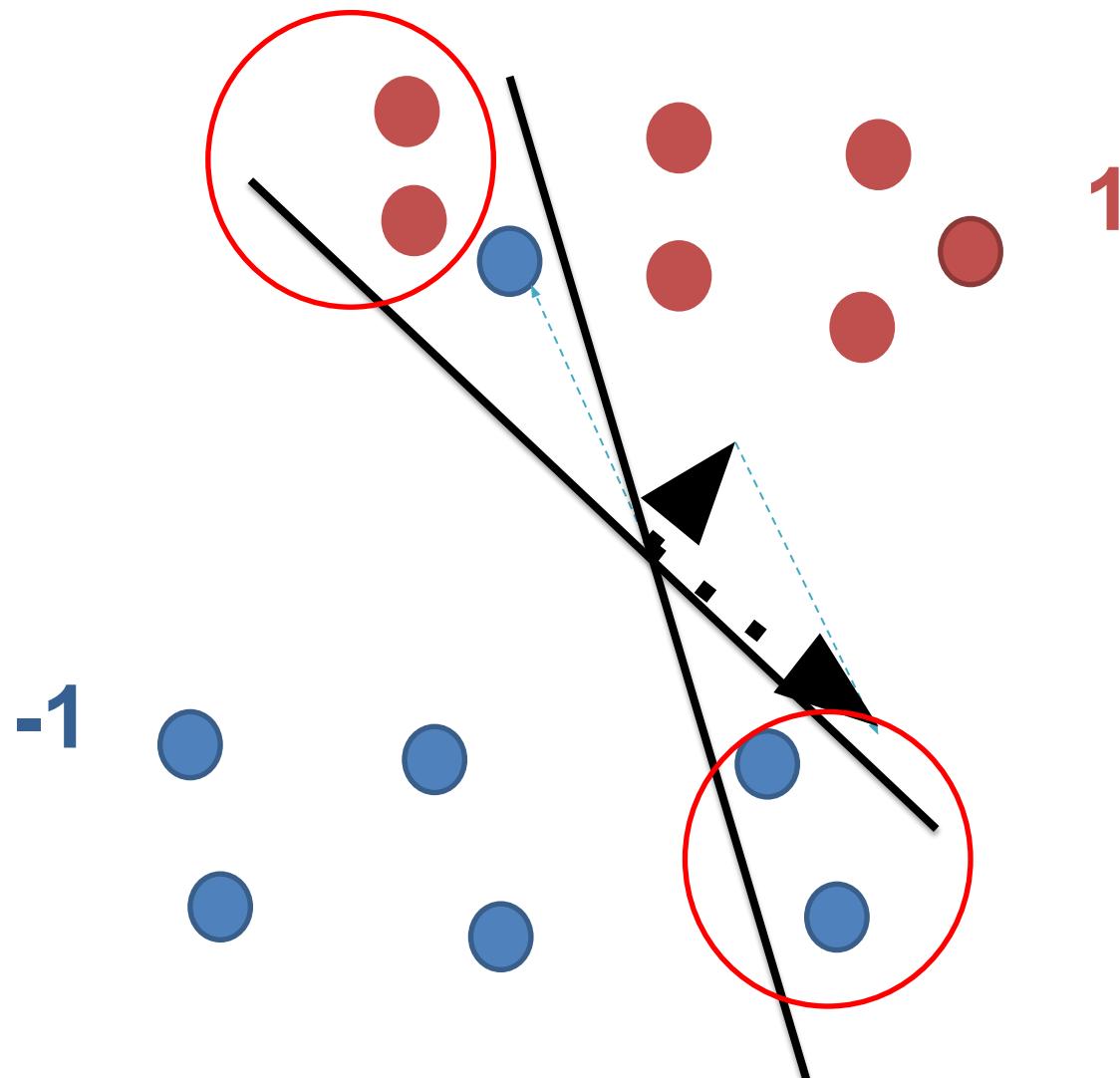
$$w^T x_i \geq 0 \quad \text{but} \quad y_i = -1$$

$$\begin{aligned} w^{t+1}^T x_i &= \left( w^t + x_i y_i \right)^T x_i \\ &= w^T x_i + y_i \|x_i\|^2 \end{aligned}$$

$$\geq 0 \quad < 0$$

Type -2

$$\begin{array}{cc} \swarrow & \downarrow \\ \leq 0 & > 0 \end{array}$$



**Fixing one error  
might lead to  
more errors elsewhere**

In general, does perceptron work for linearly separable data?

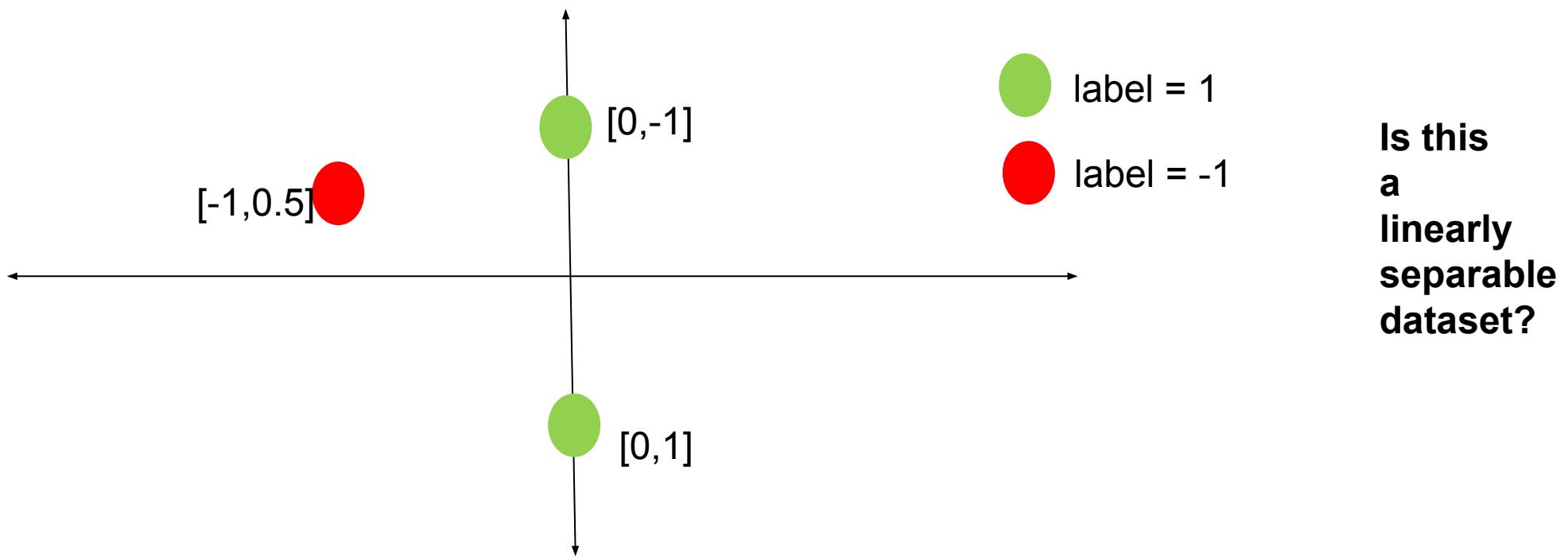
Recal

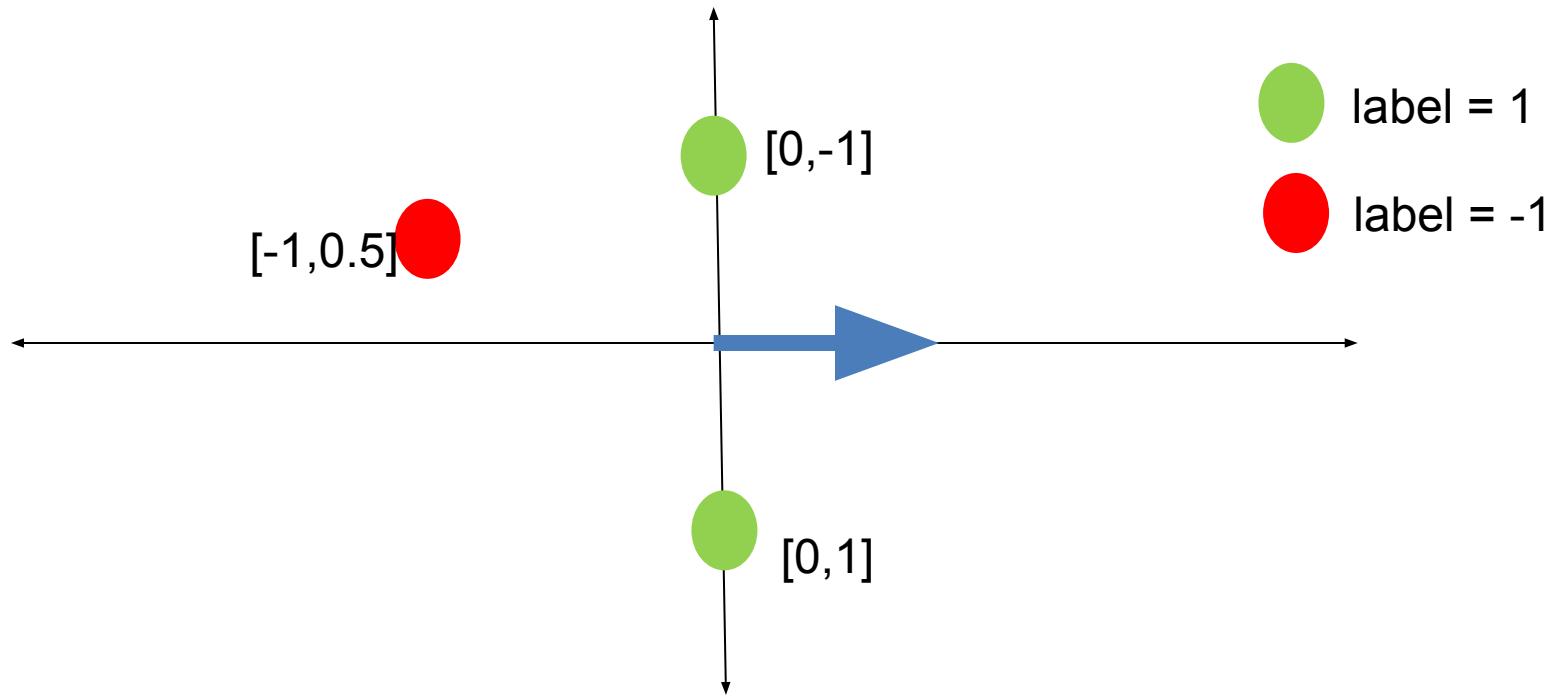
|

Lin. Sep assumption :

$$\exists w \in \mathbb{R}^d \text{ s.t } \underbrace{\text{sign}(w^\top x_i)}_{\downarrow} = y_i \quad \forall i \in [n] \\ \text{Sign}(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

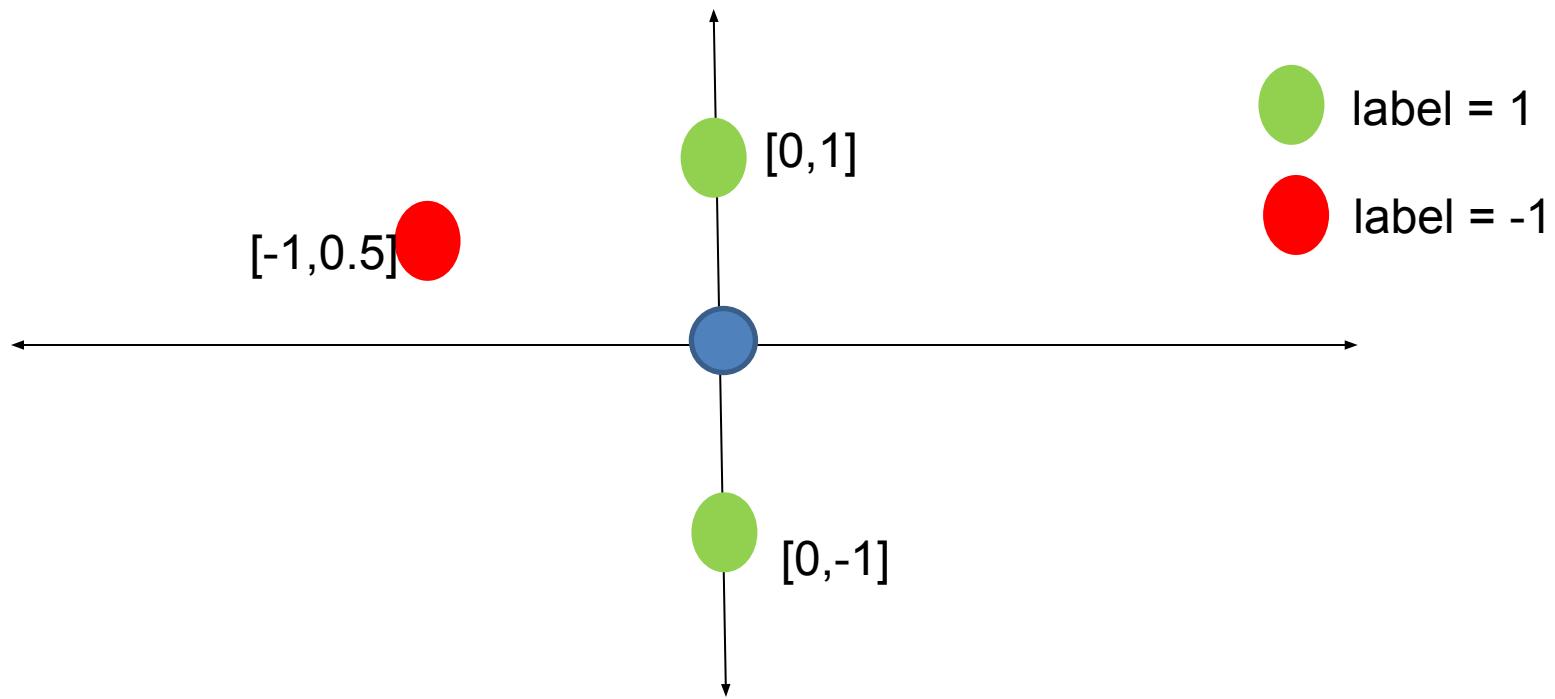
$\{1, \dots, n\}$ .





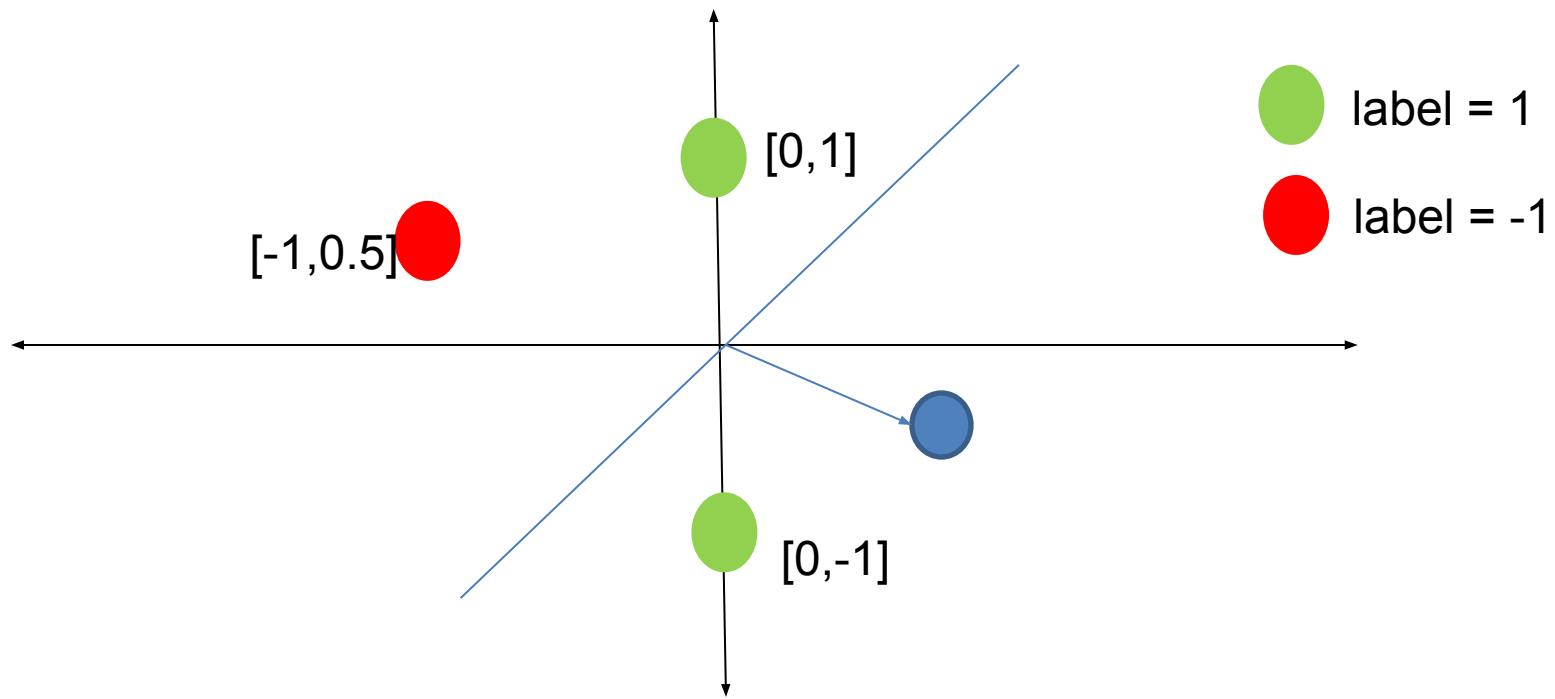
Any  $w$  in the positive x-axis linearly separates the data.  
The dataset is linearly separable.

Let's see what perceptron learns from this data!



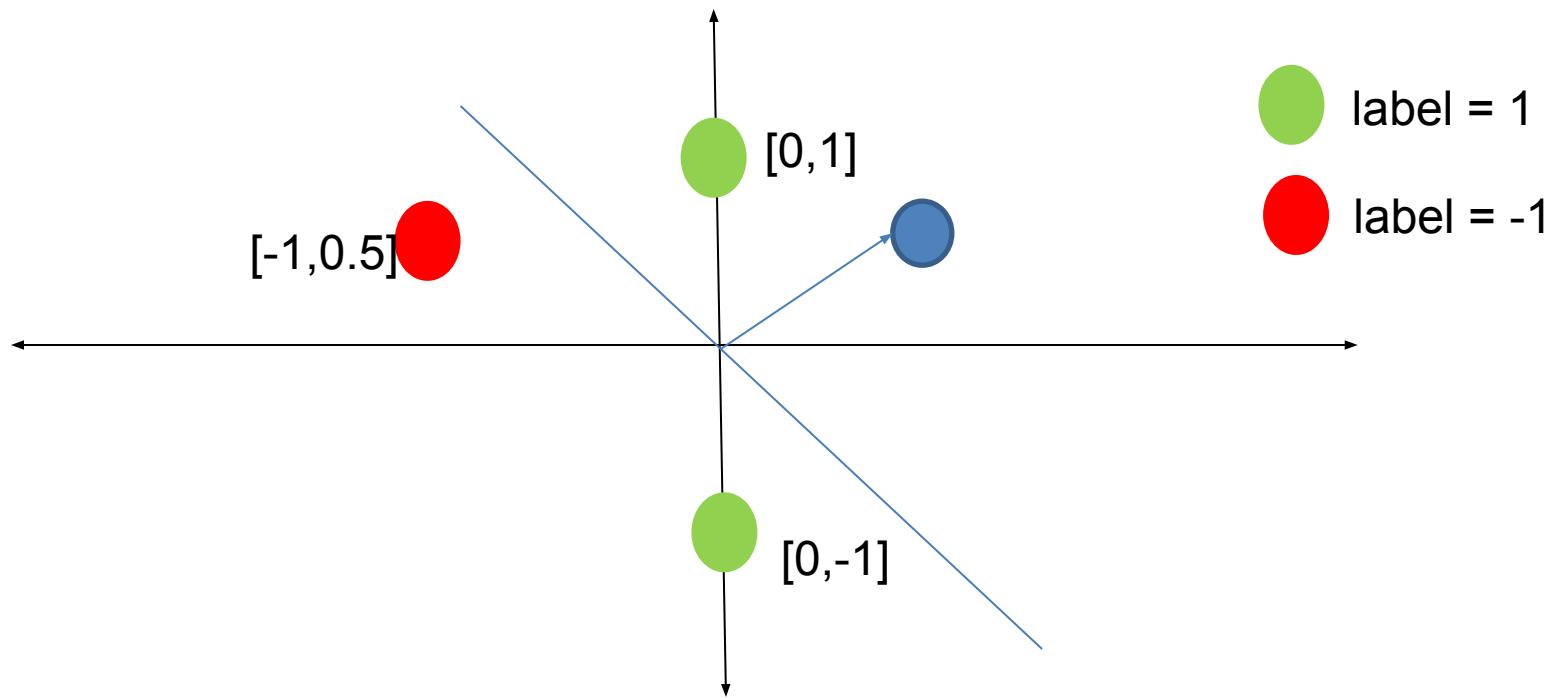
$$\mathbf{w}^0 = [0 \ 0]$$

	<b>Predicted label</b>	<b>True label</b>
$[0 \ 1]$	1	1
$[0 \ -1]$	1	1
$[-1 \ 0.5]$	1	-1



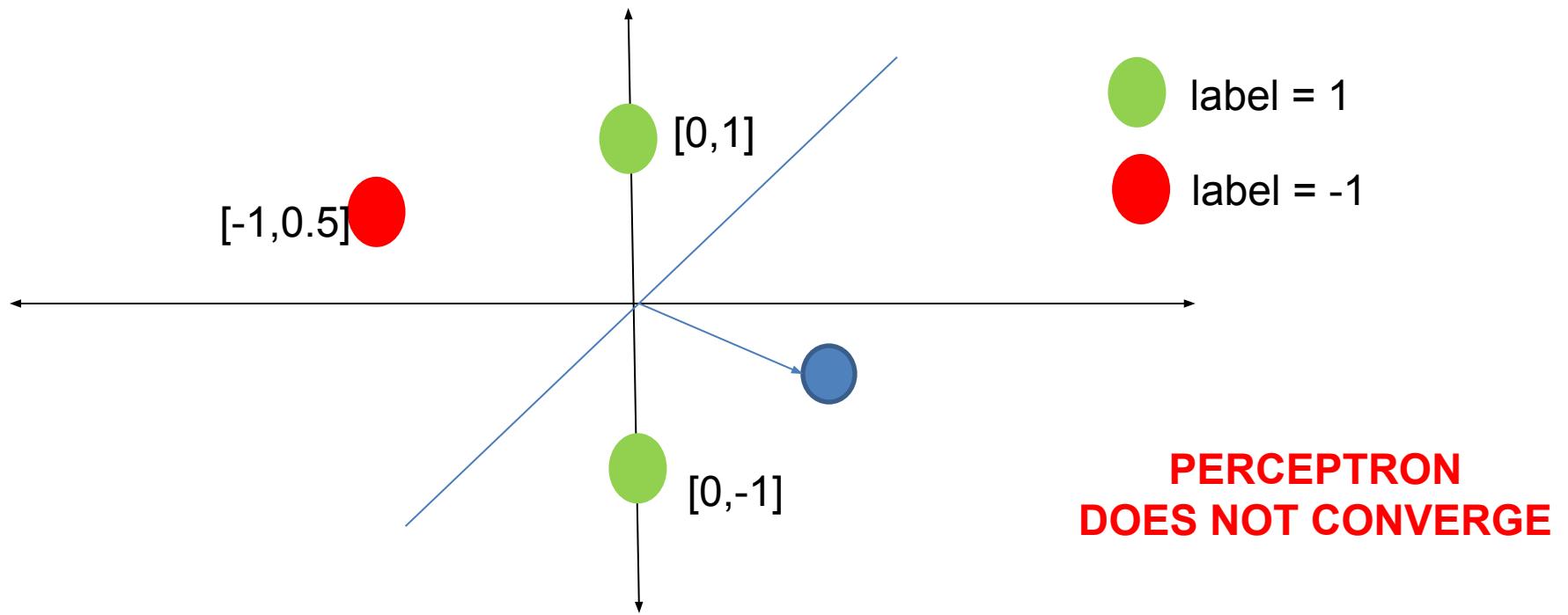
$$w^1 = [1 \ -0.5]$$

	<b>Predicted label</b>	<b>True label</b>
[0 1]	-1	1
[0 -1]	1	1
[-1 0.5]	-1	-1



$$w^2 = [1 \ 0.5]$$

	<b>Predicted label</b>	<b>True label</b>
$[0 \ 1]$	1	1
$[0 \ -1]$	-1	1
$[-1 \ 0.5]$	-1	-1



$w^3 = [1 \ -0.5]$

	Predicted label	True label
[0 1]	-1	1
[0 -1]	1	1
[-1 0.5]	-1	-1

Issue

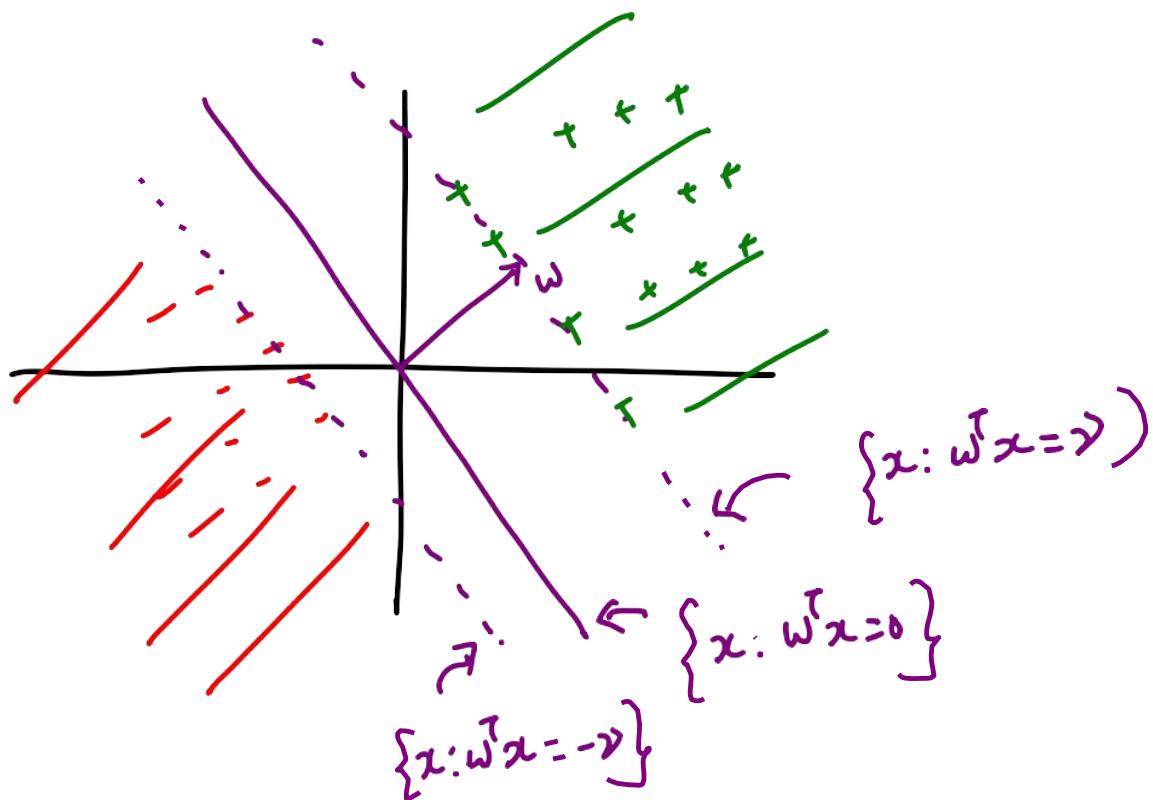
c70

- optimal  $w^* = \begin{bmatrix} c \\ 0 \end{bmatrix}$  has datapoints that lie on the linear separator.

• If we assume this isn't the case, will perception converge?

## ASSUMPTIONS

• LINEAR SEPERABILITY with  $\gamma$ -MARGIN



A Dataset  $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$  is L.S

with  $\gamma$ -margin if  $\exists w^*$  s.t

$$(w^T x_i) y_i \geq \gamma \quad \text{for some } \gamma > 0$$

The term  $w^T x_i$  is underlined in red.

# PERCEPTRON

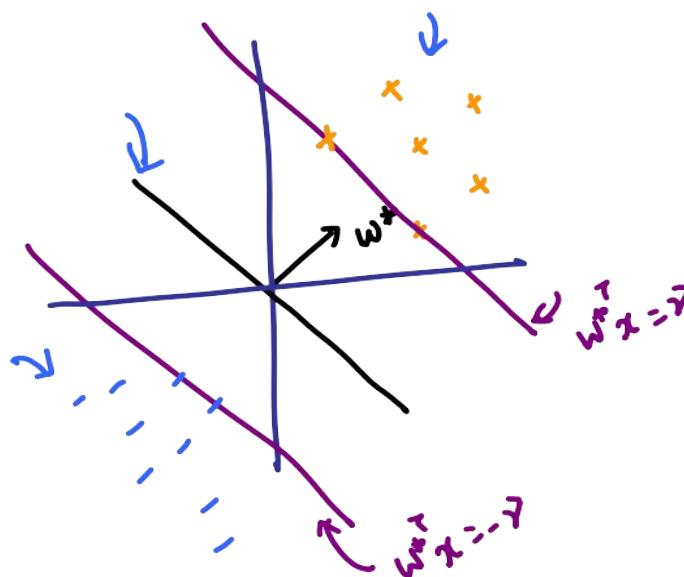
$$w^{t+1} = w^t + x_i y_i$$

## ASSUMPTIONS

- ① • Linear separability with  $\gamma$  margin

A dataset  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  is L.S with  $\gamma$  margin

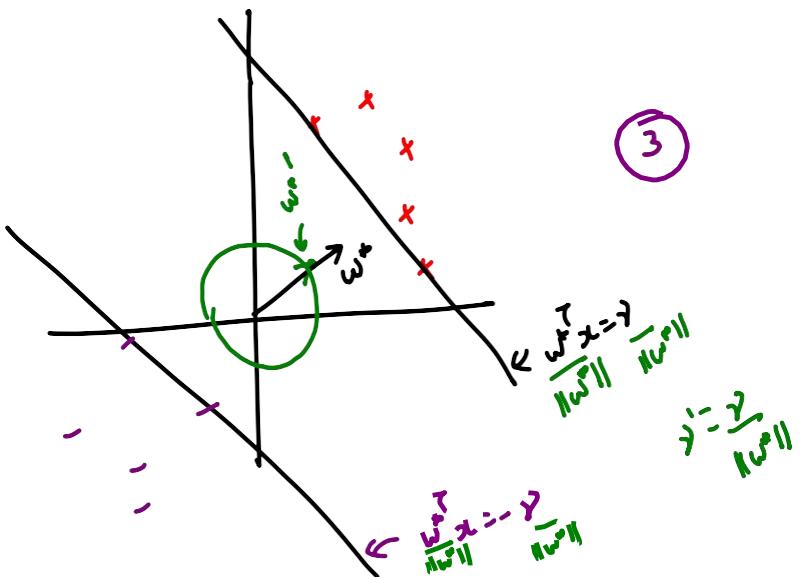
if  $\exists w^*$  s.t  $\underline{\underline{w^T x_i}} y_i \geq \gamma$   $\forall i$  for some  $\gamma > 0$



②

RADIUS    ASSUMPTION

$x_i \in D$ ,  $\|x_i\|_2 \leq R$  for some  $\underline{R > 0}$



Without loss of generality

$$\|w^*\| = 1$$

## ANALYSIS OF "mistakes" of Perceptron.

- observe that an update happens only when a "mistake" occurs.
- Say the current guess =  $w_e$  and a mistake happens w.r.t  $(x, y)$

$$w_{l+1} = w_l + x \cdot y$$

$$\|w_{l+1}\|^2 = \|w_l + x \cdot y\|^2$$

$$= (\underline{w}_l + \underline{x} \cdot \underline{y})^\top (\underline{w}_l + \underline{x} \cdot \underline{y})$$

$$\|\underline{w}_{l+1}\|^2 = \underbrace{\|\underline{w}_l\|^2}_{\leq 0} + \underbrace{2(\underline{w}_l^\top \underline{x}) \cdot \underline{y}}_{\frac{(\underline{x}^\top \underline{x}) \cdot \underline{y}^2}{\|\underline{x}\|^2}} + \underbrace{(\underline{x}^\top \underline{x}) \cdot \underline{y}^2}_{\leq R^2}$$

because  
[mistake.]

Inductively

$\|\underline{w}_{l+1}\|^2 \leq \lambda \cdot R^2$

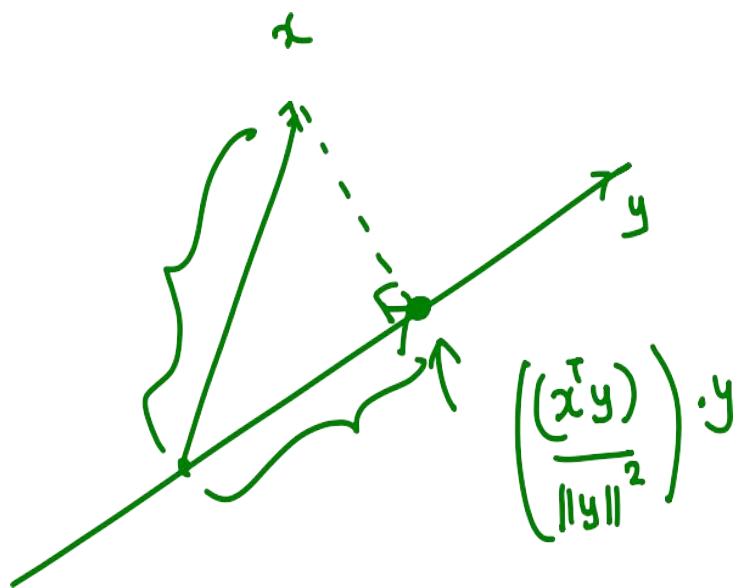
$$\leq \|\underline{w}_l\|^2 + R^2 \leq \left( \underbrace{\|\underline{w}_{l-1}\|^2 + \frac{R^2}{\lambda}}_{\|\underline{w}_0\|^2 + \lambda R^2} \right) + R^2$$

①

$$\begin{aligned}
 w_{l+1}^T w^* &= (w_l + x \cdot y)^T w^* \\
 &= w_l^T w^* + (\underbrace{w^T x}_{} \cdot y) \\
 &\geq (w_{l-1}^T w^* + \underbrace{\gamma}_{\geq \gamma}) + \gamma
 \end{aligned}$$

⋮

$w_{l+1}^T w^* \geq l \cdot \gamma$



$$\|x\|^2 \geq \left\| \left( \frac{x^T y}{\|y\|^2} \right) y \right\|^2$$

$$\geq \frac{(x^T y)^2}{\|y\|^2} \cdot \|y\|^2$$

$$\Rightarrow \underline{(x^T y)^2} \leq \underline{\|x\|^2} \underline{\|y\|^2}$$

(C.S)      Cauchy  
Schwartz      inequality       $\rightarrow$        $x^T y \leq \|x\| \|y\|$

From before

$$\underline{w_{l+1}^T w^*} \geq l \cdot \gamma$$

$$\underbrace{\|w_{l+1}\|^2}_{\text{C.S.}} \|w^*\|^2 \geq \underline{(w_{l+1}^T w^*)^2} \geq l^2 \gamma^2$$

$$\|w_{l+1}\|^2 \geq l^2 \gamma^2$$

- ②

Combining ① & ②.

$$\ell^2 \gamma^2 \leq \|w_{t+1}\|^2 \leq \ell R^2 \quad \hookrightarrow \textcircled{1}$$

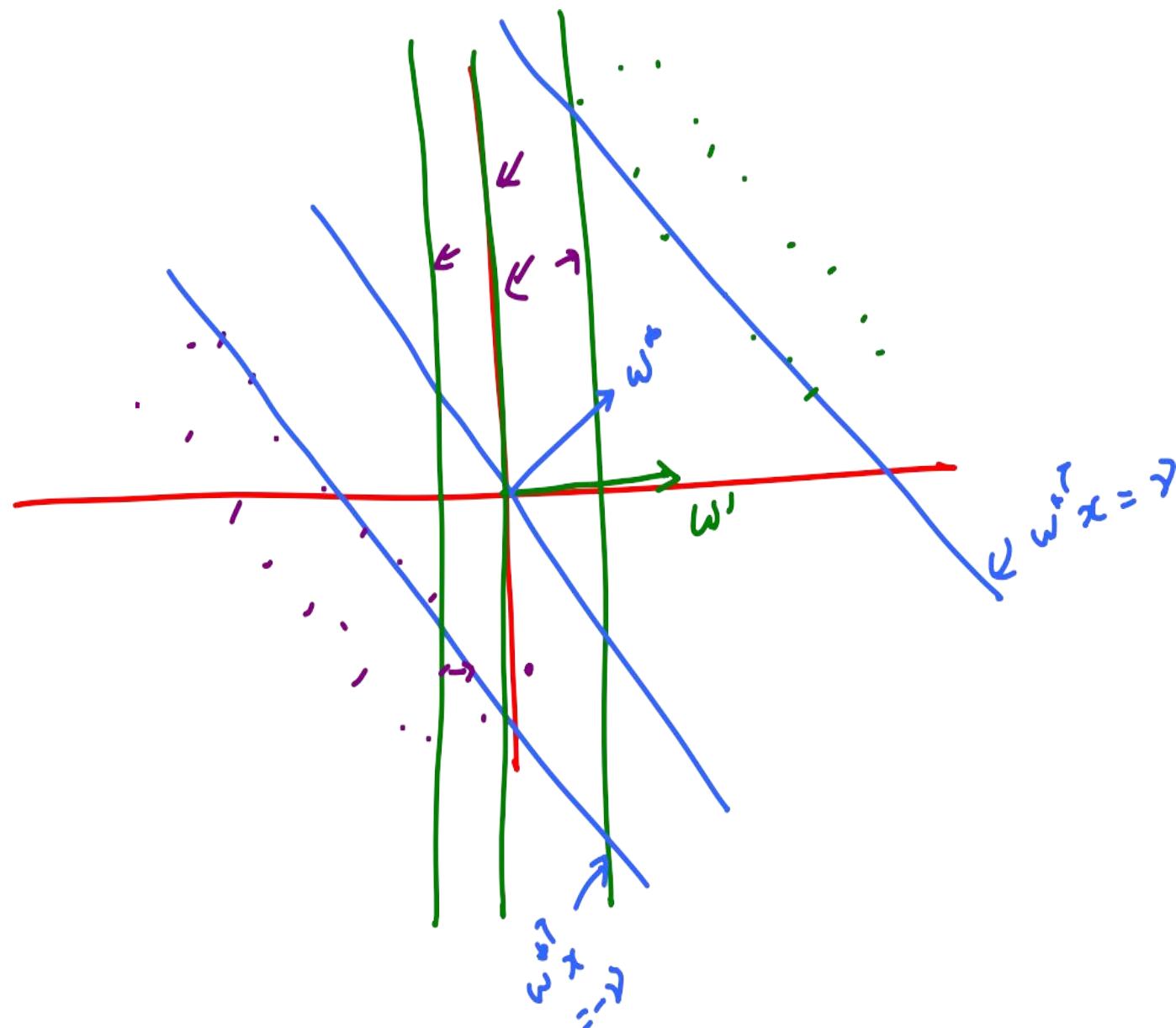
②

$$\Rightarrow \ell^2 \gamma^2 \leq \ell R^2$$

$$\boxed{\ell \leq \frac{R^2}{\gamma^2}}$$

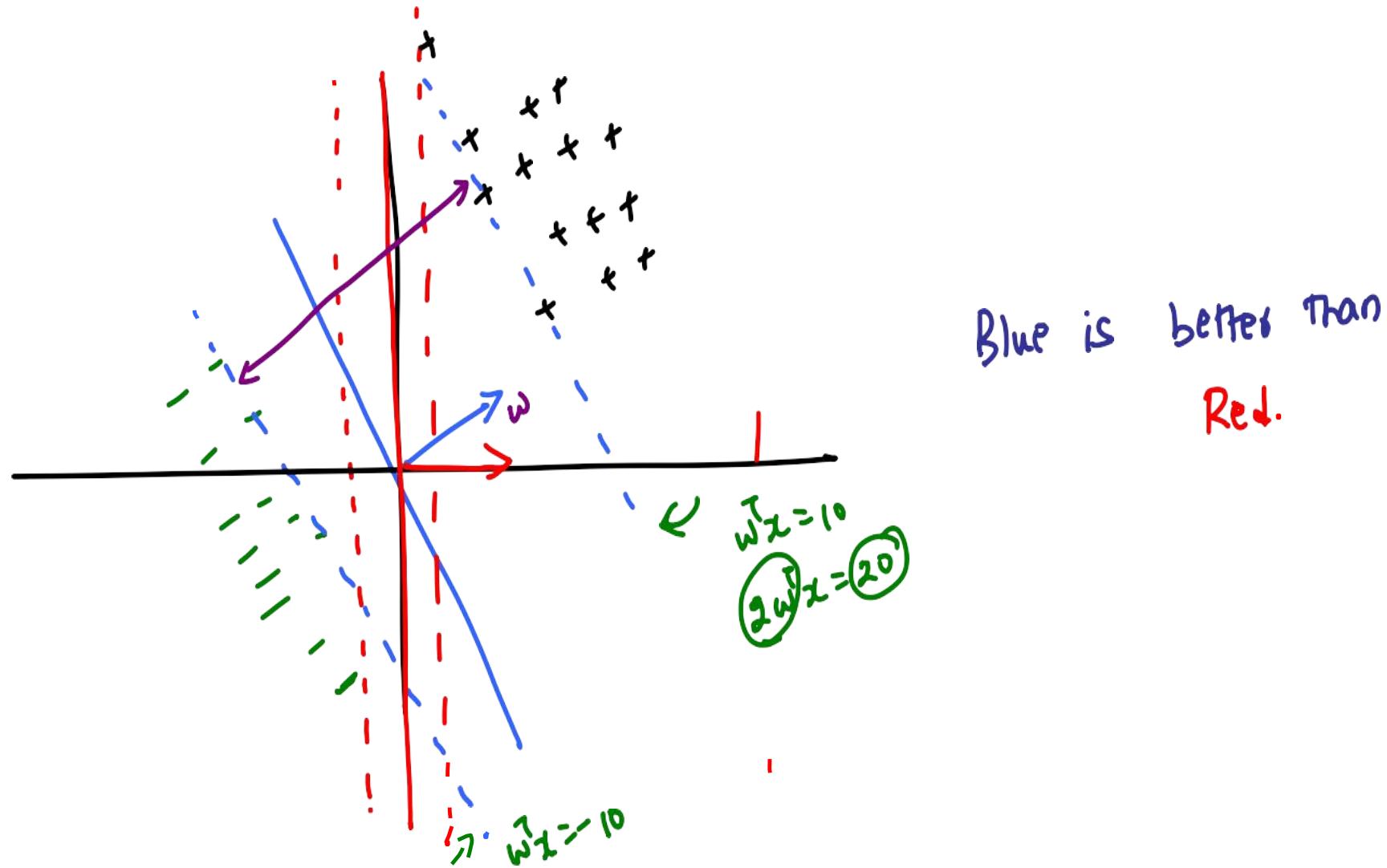
RADIUS-MARGIN  
BOUND.

$\Rightarrow$  # mistakes of Perceptron is bounded  
 $\Rightarrow$  Perceptron converged!



Perceptron's # mistakes  
depends on  $w^*$ .  
But it might  
output  $w'$ .

Goal: To come up with algorithms  
that maximize margin



Goal: Given a dataset  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

find a  $w$  with maximum width s.t. all  
points in dataset  $D$  are classified correctly

$$\max_{w, \gamma} \gamma$$

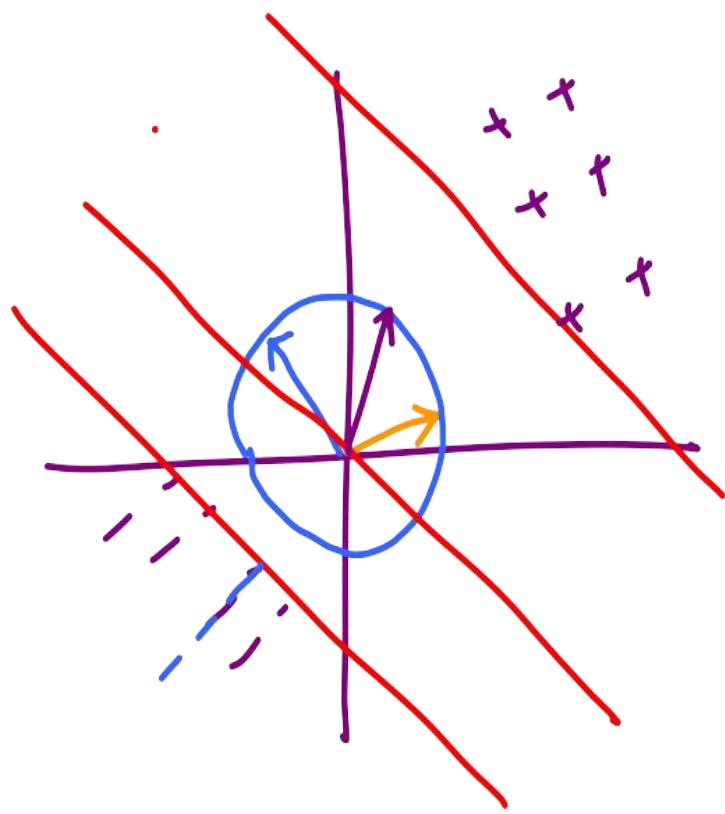
s.t.

$$(w^T x_i) y_i \geq \gamma$$

$x_i$

Issue:

Can scale  $w$   
arbitrarily

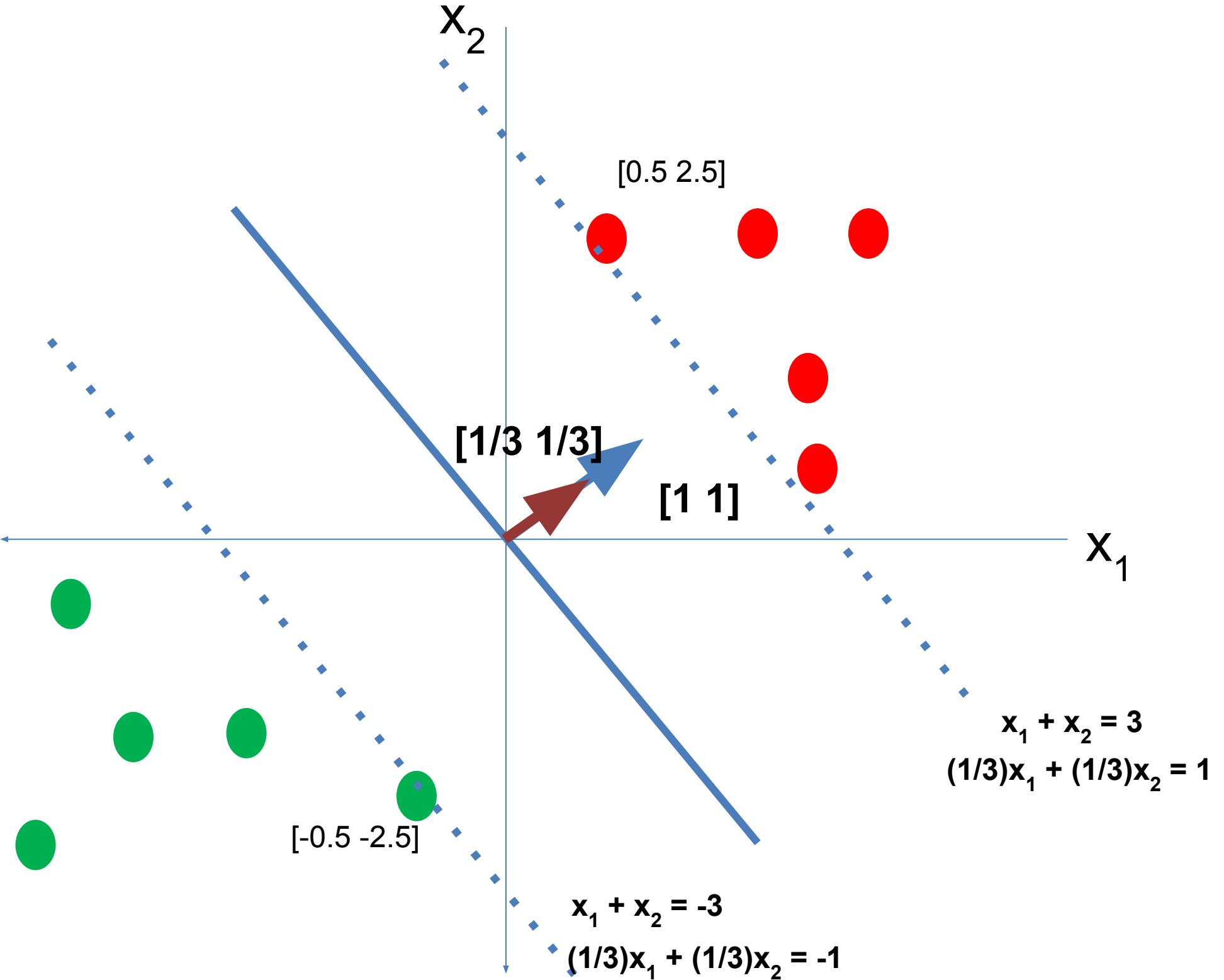


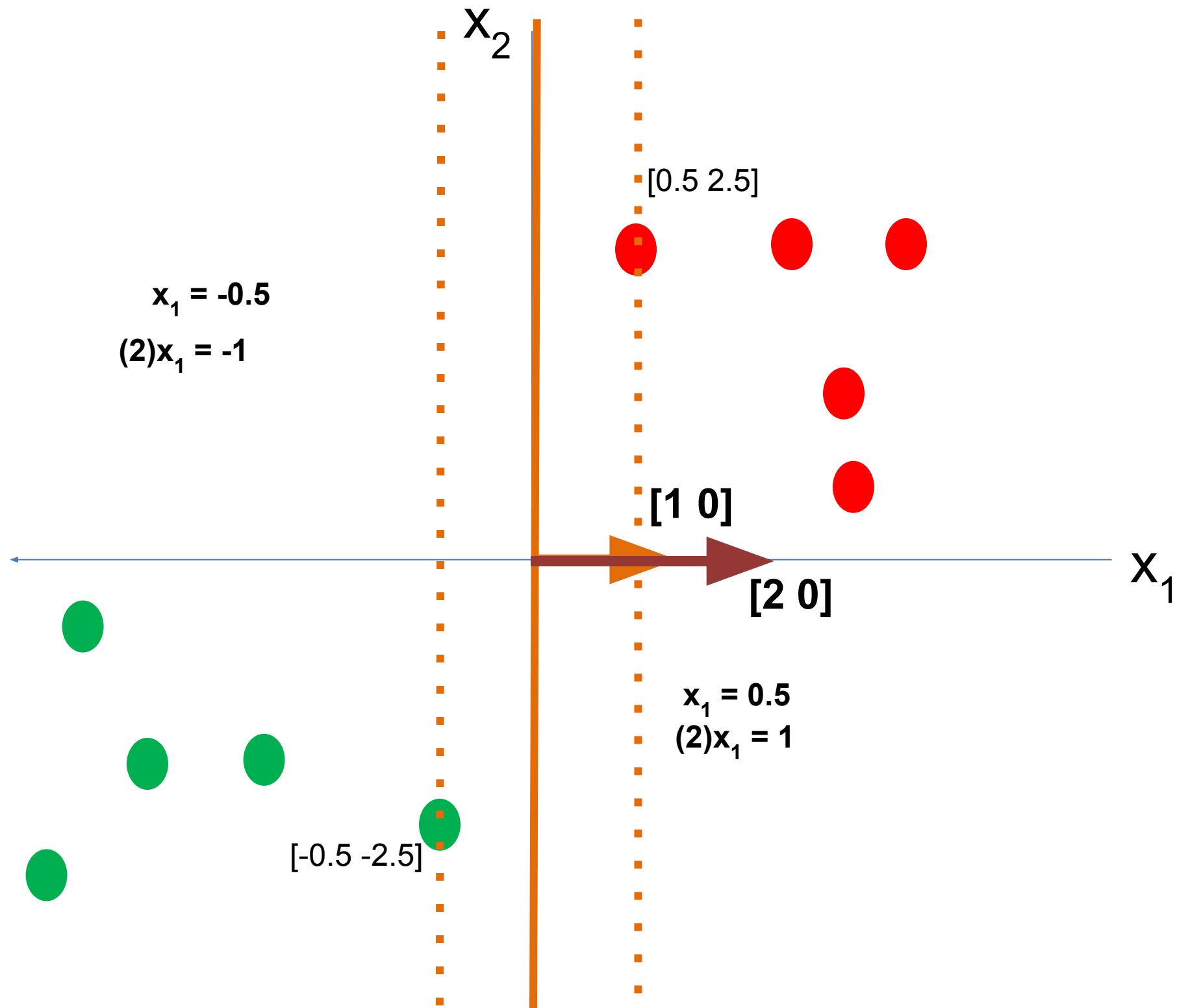
Possible fix

$$\max_{w, \gamma} \gamma$$

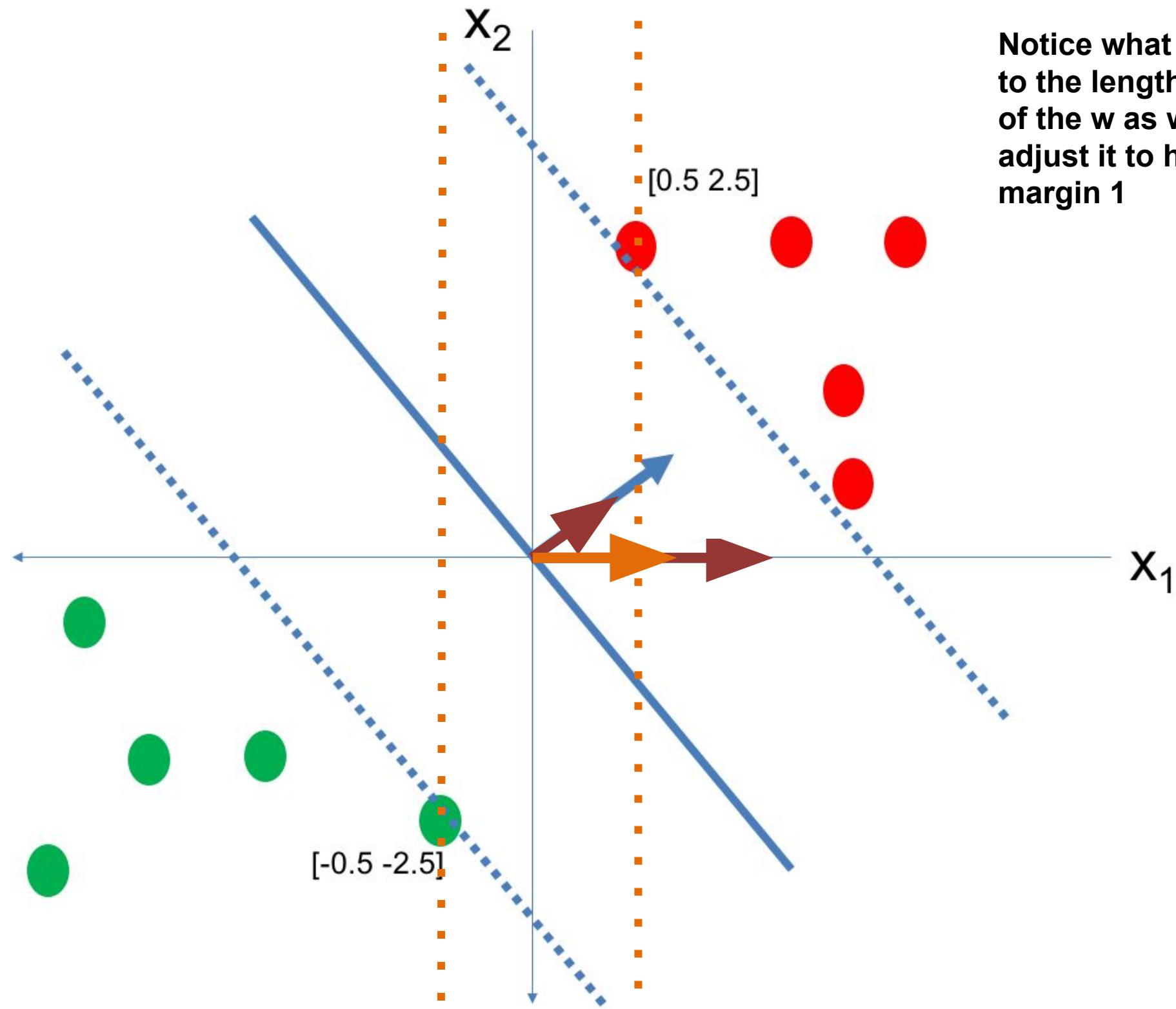
$$(w^T x_i) y_i \geq \gamma$$

$$\|w\|^2 = 1$$





Notice what happens  
to the lengths  
of the  $w$  as we  
adjust it to have  
margin 1

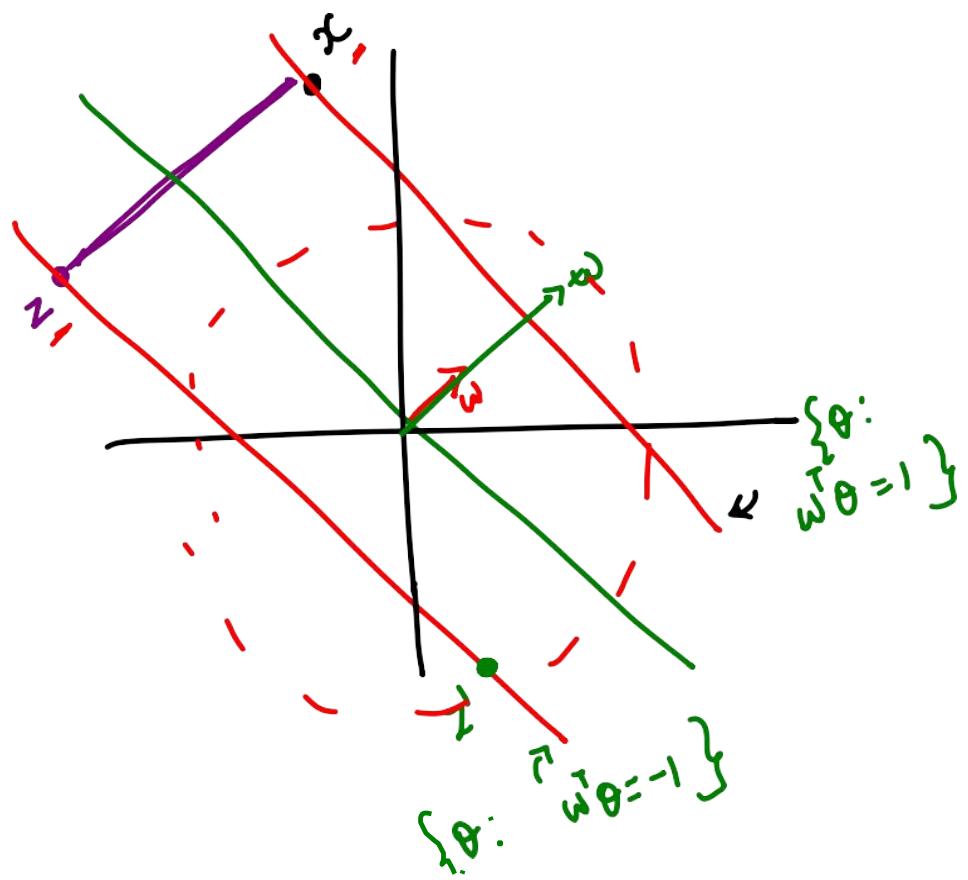


## OBSERVATIONS

- > Once a direction is fixed, the width between the margin lines is fixed
- > If the width is large, then the  $w$  that achieves margin 1 in that direction has smaller length
- > If the width is small, then the  $w$  that achieves margin 1 in that direction has larger length
- > In general,  $\text{width}(w)$  seems to be inversely proportional to  $\|w\|$

$$\begin{array}{l} \max_w \\ \text{width}(w) \\ \text{s.t. } (\mathbf{w}^\top \mathbf{x}_i) y_i \geq 1 \end{array}$$

What is  $\text{width}(w)$ ?



$$\begin{aligned} \min_{z:} \quad & \frac{1}{2} \|x - z\|^2 \\ \text{s.t.} \quad & w^T x = 1 \\ & w^T z = -1 \end{aligned}$$

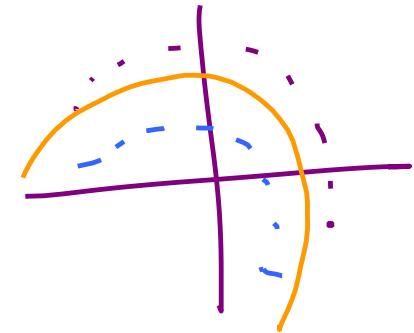
[Exercise]

$$\text{width}(w) = \frac{2}{\|w\|^2}$$

$$\max_w$$

$$\frac{1}{2} \|\omega\|^2$$

$$\text{s.t. } \forall i \quad (\omega^\top x_i) y_i \geq 1$$



$$\min_w$$

$$\frac{1}{2} \|\omega\|^2$$

$$\text{s.t. } \forall i \quad (\omega^\top x_i) y_i \geq 1$$

### Issues

- L.S is a strong assumption

- Non-linear structure?

DETOUR

$$\begin{array}{l} \min_{\omega} f(\omega) \\ g(\omega) \leq 0 \end{array}$$



$$\underline{L}(\omega, \alpha) = f(\omega) + \alpha \cdot g(\omega)$$

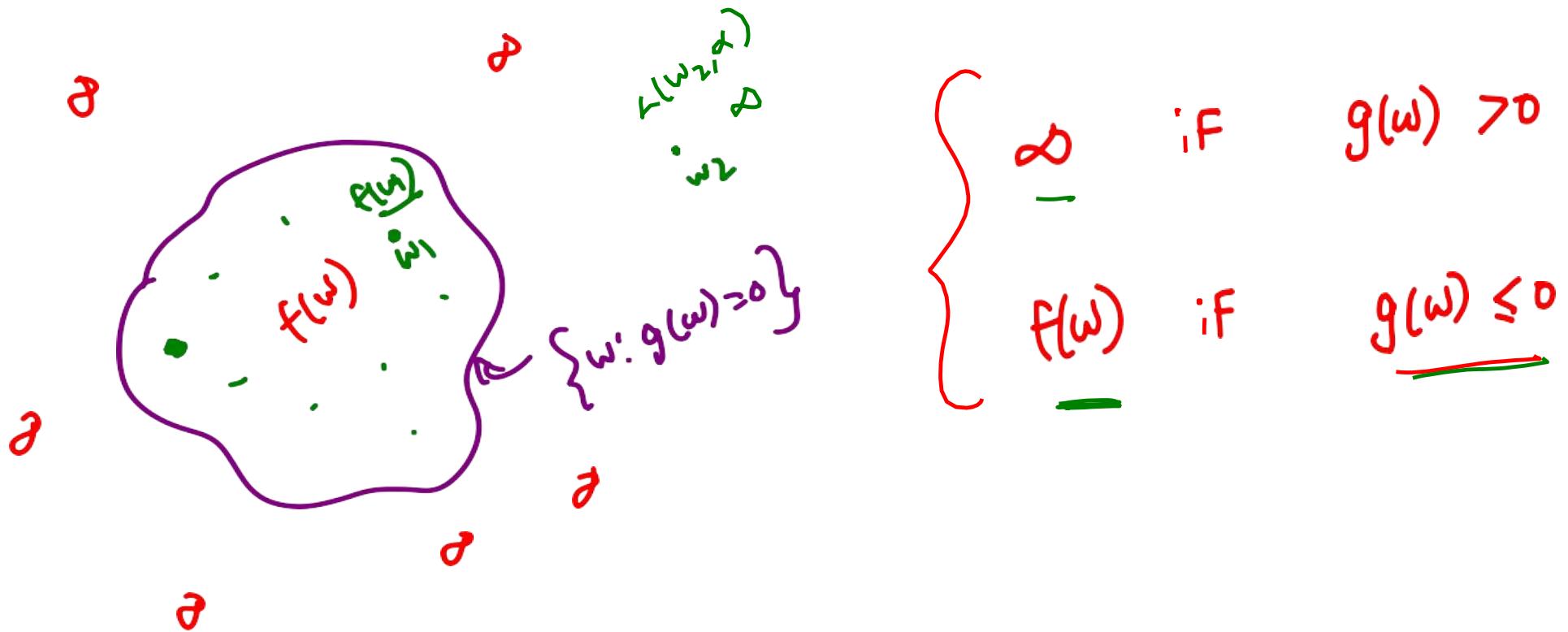
Fix some  $\omega$ .

Consider

$$\max L(w, \alpha)$$

$$\alpha \geq 0$$

$$= \max_{\alpha \geq 0} \underline{f(w)} + \alpha \underline{g(w)}$$



$$\min_w \left[ \max_{\alpha \geq 0} \frac{R(w)}{\lambda(w, \alpha)} \right]$$

$\xrightarrow{\text{equivalent}}$

$$\begin{aligned} & \min_w f(w) \\ & \text{s.t. } g(w) \leq 0 \end{aligned}$$

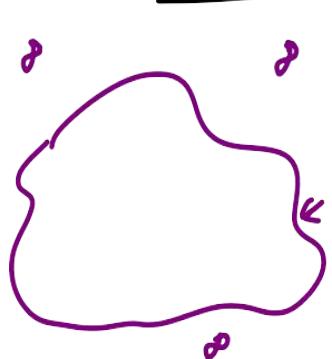
}

## SO FAR

MAX-Margin formulation

$$\begin{array}{ll} \min_{w \in \mathbb{R}^d} & \frac{1}{2} \|w\|^2 \\ \text{s.t.} & (w^\top x_i) y_i \geq 1 \end{array}$$

- A



DETOUR

$$\begin{array}{ll} \min_w & f(w) \\ \text{s.t.} & g(w) \leq 0 \end{array}$$

$$\min_w \left[ \max_{\alpha \geq 0} [f(w) + \alpha g(w)] \right]$$

- Can we swap min and max?

## Multiple Constraints

→ Same idea

$$\min_{\omega} f(\omega)$$

$$\text{s.t. } g_i(\omega) \leq 0 \quad \begin{matrix} +i \\ =1..K \end{matrix}$$

$$\min_{\omega} \left[ \max_{\substack{\{\alpha_1, \dots \\ \alpha_K \geq 0\}}} \left[ f(\omega) + \underbrace{\sum_{i=1}^k \alpha_i g_i(\omega)}_{\text{---}} \right] \right]$$

III Strong duality for convex  $f, g_i$

$$\max_{\substack{\alpha_1, \dots, \alpha_K \geq 0}} \min_{\omega} f(\omega) + \sum_{i=1}^k \alpha_i g_i(\omega)$$

$$\min_{\omega} \frac{1}{2} \|\omega\|^2 \quad \leftarrow f(\omega)$$

s.t.  $(\omega^T x_i) y_i \geq 1$

$\underbrace{(\omega^T x_i) y_i}_{+i} \geq 1$

$1 - (\omega^T x_i) y_i \leq 0$

$$g_i(\omega) = 1 - (\omega^T x_i) y_i$$

$$L(\omega, \alpha) = \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \alpha_i (1 - (\omega^T x_i) y_i)$$

$\uparrow$   
 $\in \mathbb{R}^n$

$$\min_{\omega} \left[ \max_{\alpha \geq 0} \left[ \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \alpha_i (1 - (\omega^T x_i) y_i) \right] \right]$$

$\hookrightarrow \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{bmatrix} \geq 0$

|||

$$\max_{\alpha \geq 0} \left[ \min_{\omega} \left[ \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \alpha_i (1 - (\omega^T x_i) y_i) \right] \right]$$

Fix  $\alpha \geq 0$

$$\min_w \left[ \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i (1 - w^T x_i) y_i \right]$$

Grad w.r.t  $w$

$$w^* + \sum_{i=1}^n -\alpha_i x_i y_i = 0$$

$$w^* = \sum_{i=1}^n \alpha_i x_i y_i$$

$\in \mathbb{R}^d$   
 $\{\alpha_i^{+1, -1}\}$   
 Fixed  
 Choice

In matrix notation

$$w^* = X Y \alpha$$

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{bmatrix}_{d \times n} \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}_{n \times 1}$$

Substituting  $\underline{\text{soln}}$  back in the objective.

$$\begin{aligned} & \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^n \alpha_i (1 - \omega^\top x_i) y_i \\ &= \frac{1}{2} \underline{\omega^\top \omega} + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i (\omega^\top x_i) y_i \end{aligned}$$

$$\frac{1}{2} (\mathbf{x} \mathbf{y} \boldsymbol{\alpha})^T (\mathbf{x} \mathbf{y} \boldsymbol{\alpha}) + \boldsymbol{\alpha}^T \mathbf{1} - \sum_{i=1}^n (\mathbf{x} \mathbf{y} \boldsymbol{\alpha})^T \mathbf{x}_i y_i \boldsymbol{\alpha}_i$$

$\begin{bmatrix} \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_n \end{bmatrix} \quad \begin{bmatrix} \cdot \\ \vdots \\ \cdot \end{bmatrix}$

on Simplification [please-do this]

$$\boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} (\mathbf{x} \mathbf{y} \boldsymbol{\alpha})^T (\mathbf{x} \mathbf{y} \boldsymbol{\alpha})$$

### DUAL PROBLEM

Giving  $\mathbf{y}$  instead of  $\mathbf{x}$

$$\max_{\boldsymbol{\alpha} \geq 0}$$

$\in$  easy constraints

$$\boldsymbol{\alpha}^T \mathbf{1} - \frac{1}{2} \boldsymbol{\alpha}^T \mathbf{y}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\alpha}$$

can be  
KERNELIZED!

$n \times n$ .

## Revisiting The Lagrangian

$$\min_w \left[ \max_{\alpha > 0} f(w) + \alpha g(w) \right] = \max_{\alpha > 0} \left[ \min_w f(w) + \alpha g(w) \right]$$

PRIMAL

DUAL

$$w^* \quad \quad \quad \alpha^*$$
$$\boxed{\max_{\alpha > 0} f(w^*) + \alpha^* g(w^*)} = \min_w f(w) + \alpha^* g(w)$$

$$f(\omega^*) = f(\omega) + \alpha^* g(\omega')$$

$$f(\omega^*) \leq f(\omega^*) + \alpha^* g(\omega^*)$$

$$\Rightarrow \boxed{\alpha^* g(\omega^*) \geq 0}$$

But we know  $\alpha^* g(\omega^*) \leq 0$

$\Rightarrow \boxed{\alpha^* g(\omega^*) = 0} \rightarrow$  COMPLEMENTARY SLACKNESS

For multiple constraints

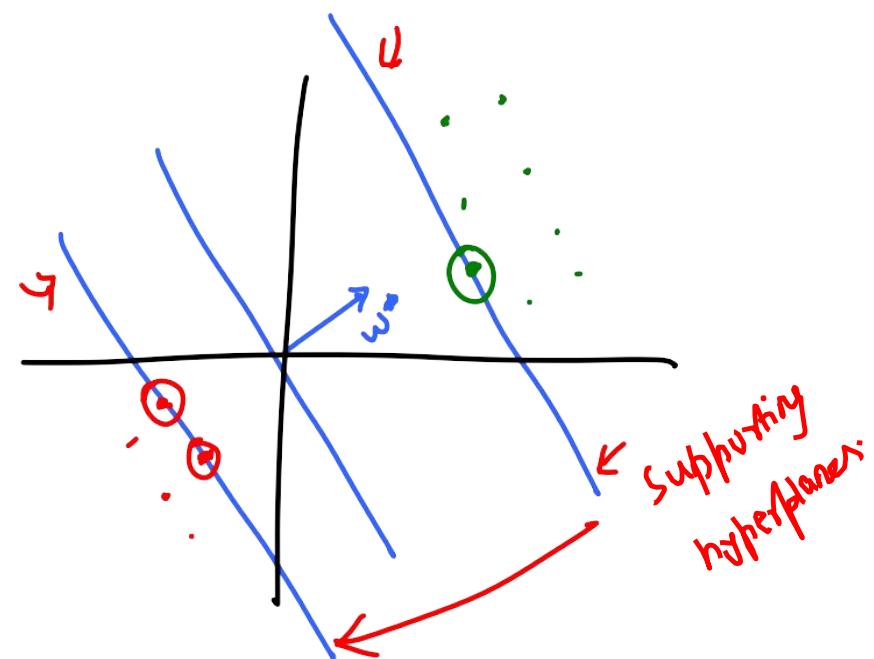
For our problem

$$\alpha_i^* g_i(\omega^*) = 0 \quad \forall i$$

$$(\alpha_i^*) (1 - (\omega^T x_i) y_i) = 0 \quad \forall i$$

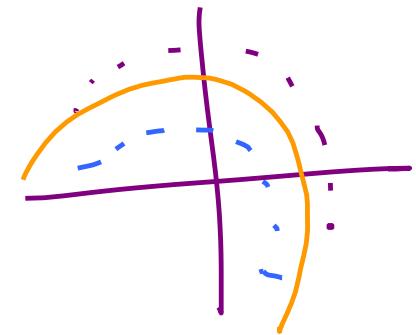
$$\omega^* = \sum_{i=1}^n \alpha_i^* x_i y_i$$

If  $\alpha_i^* > 0 \Rightarrow 1 - (\omega^T x_i) y_i = 0$   
 $\Rightarrow (\omega^T x_i) y_i = 1$



$$\min_{\omega} \frac{1}{2} \|\omega\|^2$$

$$\text{s.t. } \forall i \quad (\omega^\top x_i) y_i \geq 1$$



### Issues

- L.S is a strong assumption

- Non-linear structure?