

PRIMAL WAS  $\theta P$ , DUAL IS  $\theta P$ ,

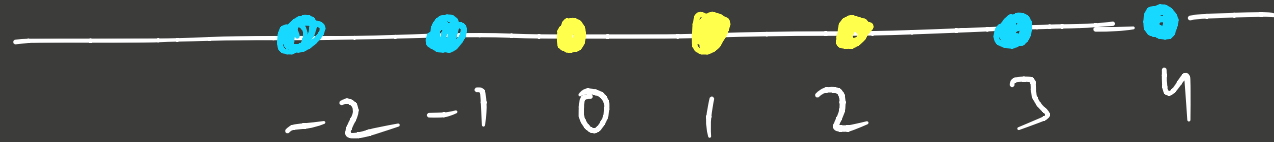
WHY BOTHER CONVERTING TO DUAL?

JUST WAIT FOR FEW MORE MINS

ANSWER: "KERNEL TRICK"

# NON-LINEARLY SEPARABLE DATA

---

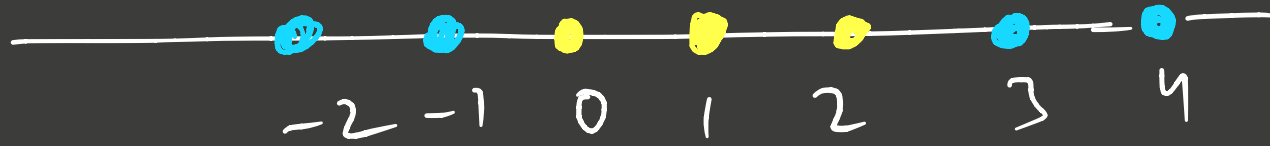


Data not separable in  $\mathbb{R}$

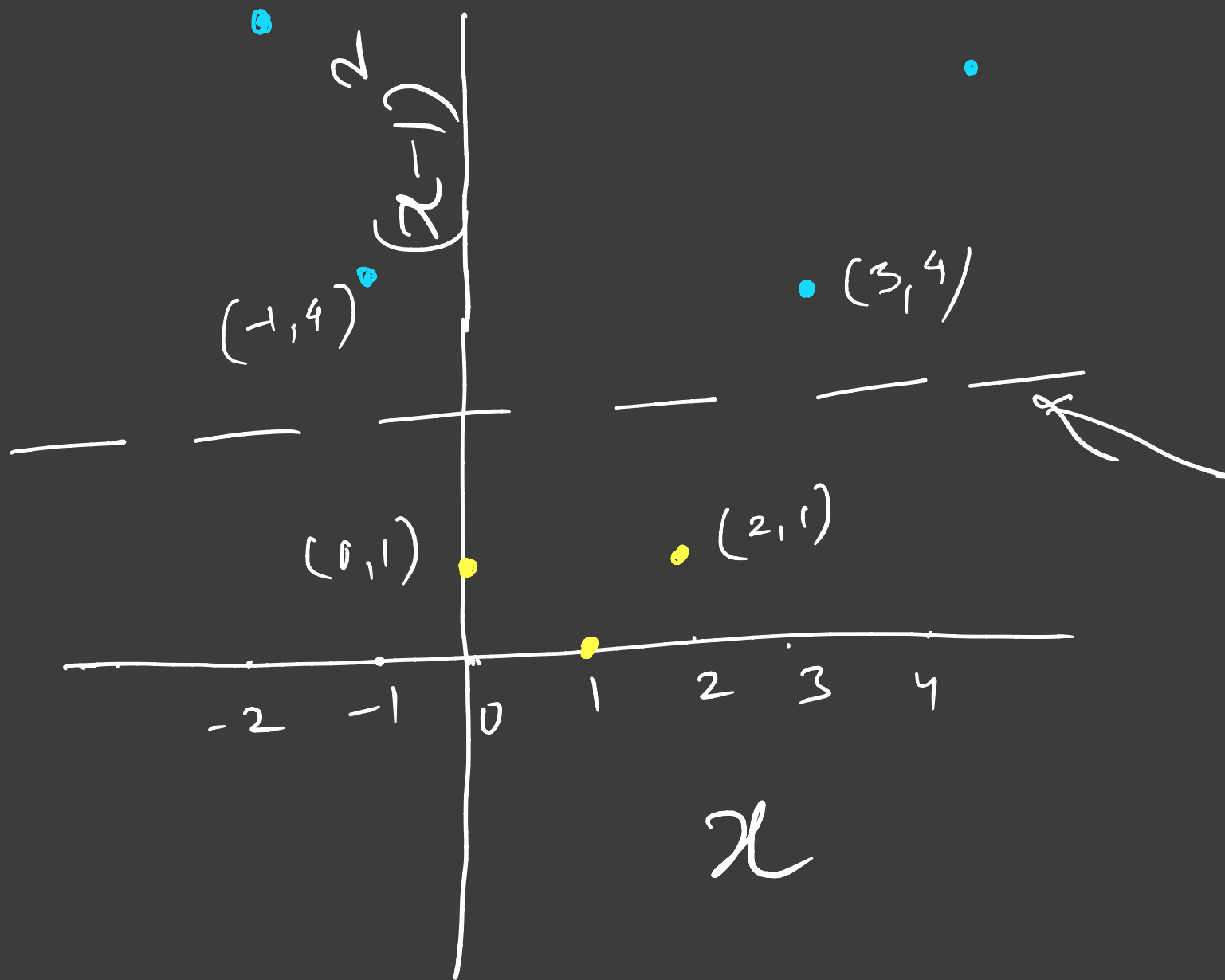
Can we still use SVM?

Yes!

How: Project data to a higher dimensional space.

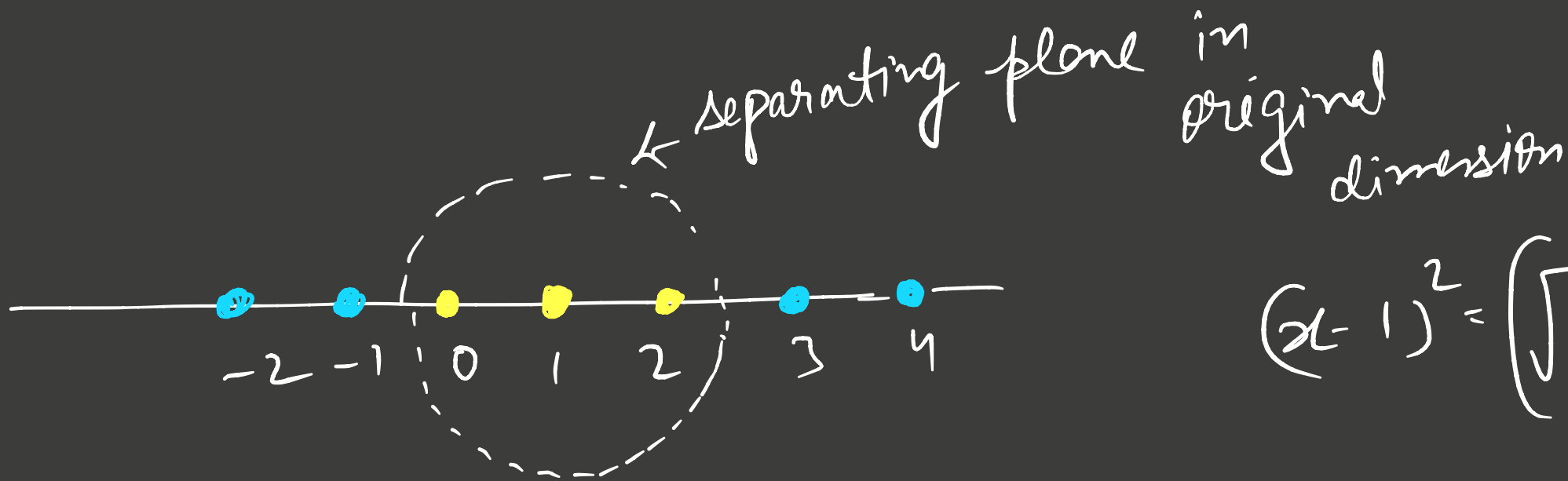


original data  
in  $\mathbb{R}$



Transformed  
data in  
 $\mathbb{R}^2$

max margin  
separating  
hyperplane



$$(x-1)^2 = \left(\sqrt{\frac{5}{2}}\right)^2$$

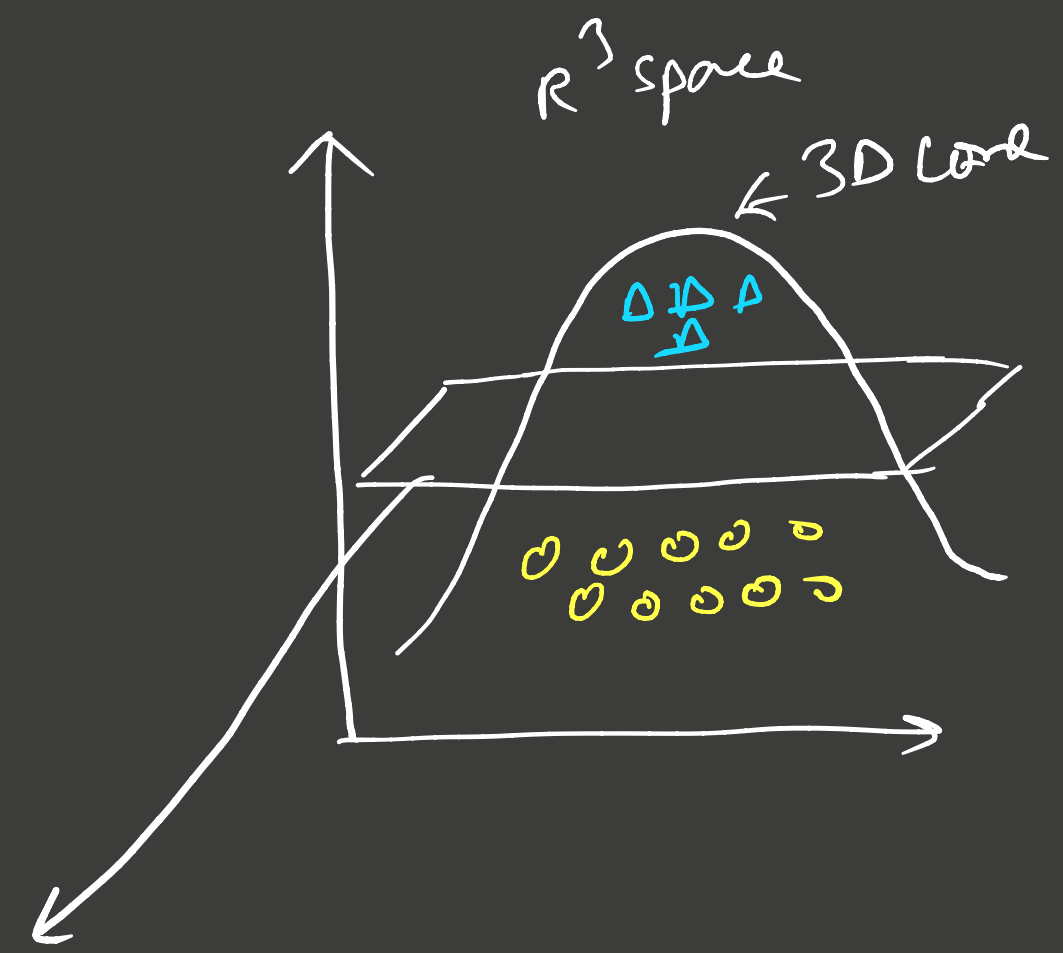
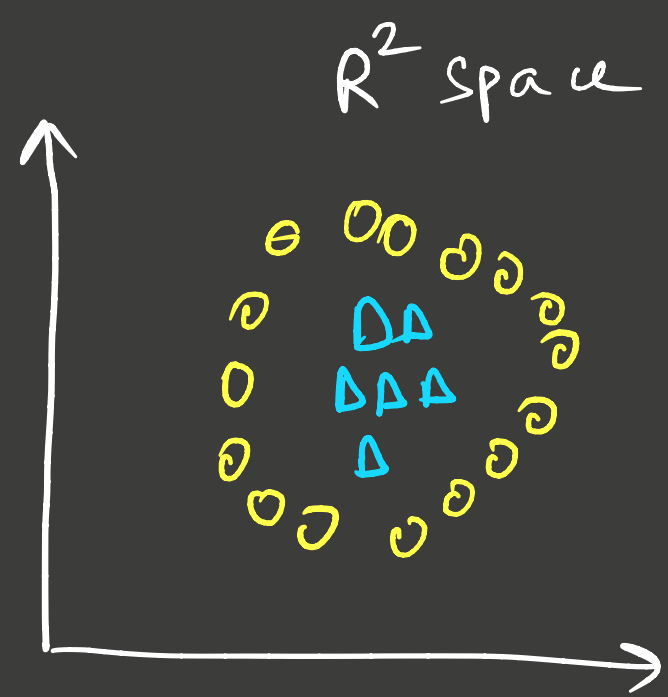
Circle :

Center ( $x=1$ )

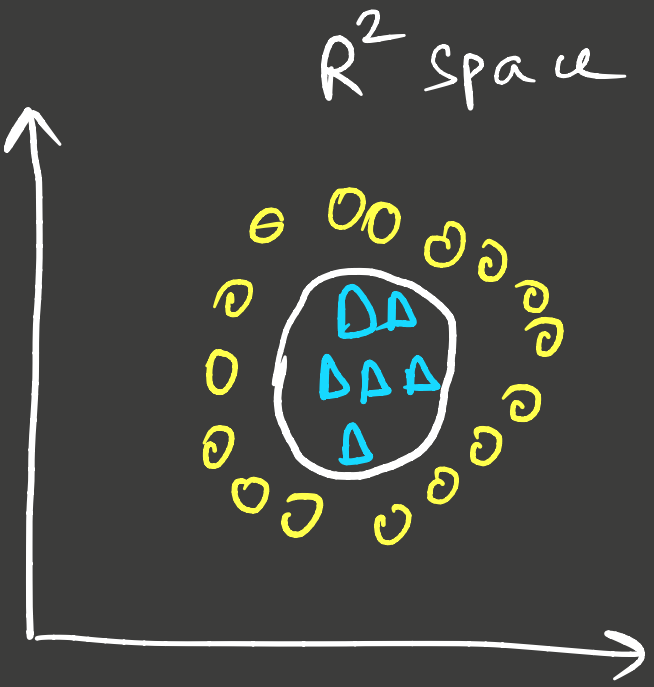
Radius  $\sqrt{5/2}$



# ANOTHER EXAMPLE TRANSFORMATION



Equivalent in  $\mathbb{R}^2$  is



# Projection / Transformation Function

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$$

where  $d =$  original dimension

$D =$  New dimension

In our example;

$$d = 1; D = 2$$

Linear SVM

MAXIMIZE

$$L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \vec{x}_i \cdot \vec{x}_j$$

s.t.

CONSTRAINTS

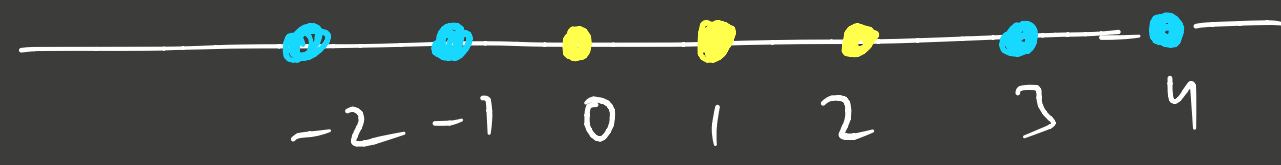


TRANSFORMATION( $\phi$ )



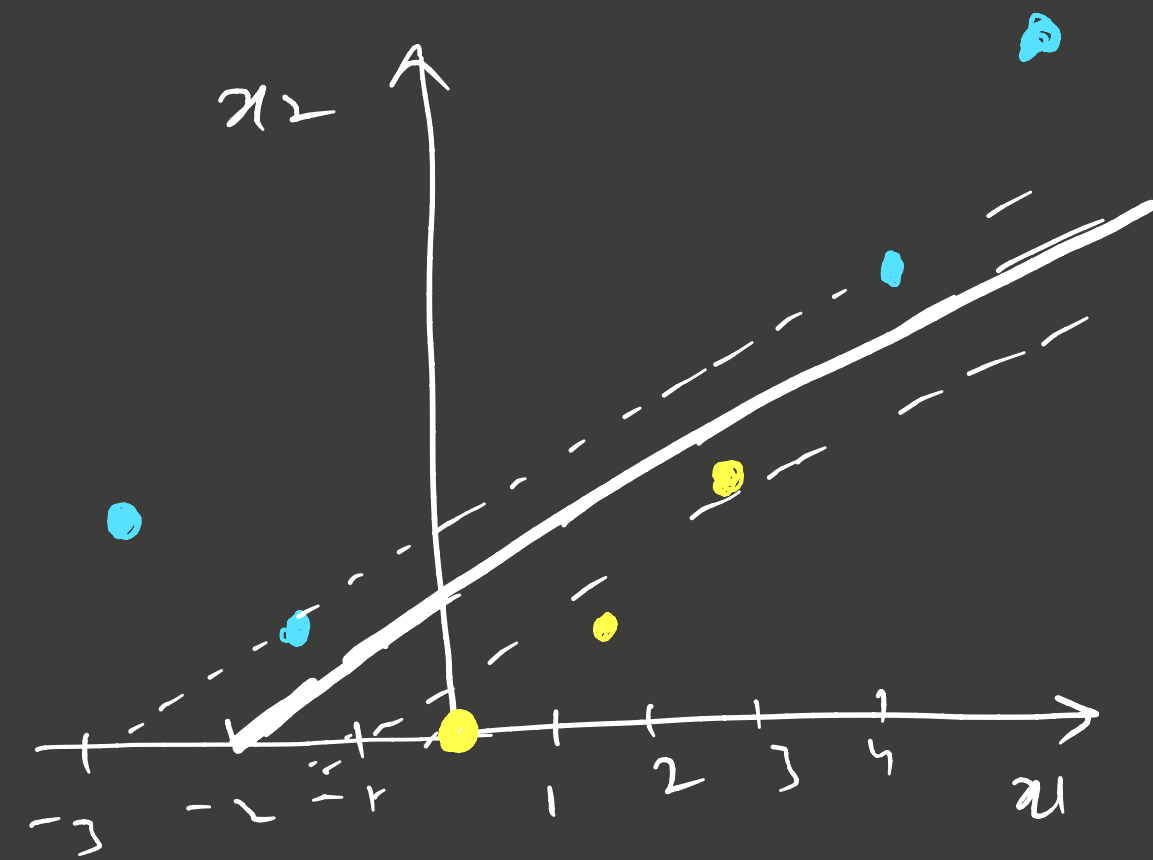
$$L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

# TRIVIAL EXAMPLE (again)



Original Data  $(x) \in \mathbb{R}$

Transformed Data  $(\phi(x) = \langle \sqrt{2}x, x^2 \rangle)$



## Steps

① Compute  $\phi(x)$  for each point

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$$

② Compute dot products over  $\mathbb{R}^D$  space

③ If  $D \gg d$

Both steps are expensive!

# KERNEL TRICK

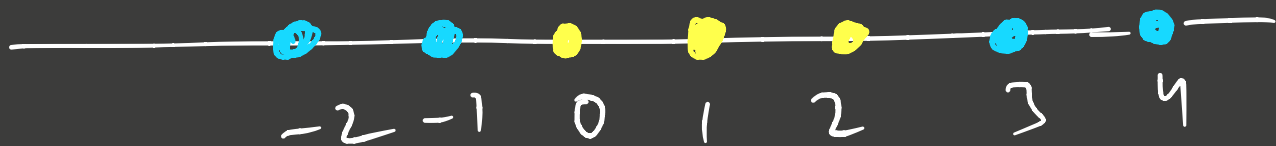
Can we compute  $K(\bar{x}_i, \bar{x}_j)$

s.t.

$$K(\bar{x}_i, \bar{x}_j) = \phi(\bar{x}_i) \cdot \phi(\bar{x}_j)$$

Some func<sup>n</sup> of  
dot product in  
original dimension

Dot product in high  
dimensions (after  
transformation)



$$\phi(x) = \langle \sqrt{2}x, x^2 \rangle$$

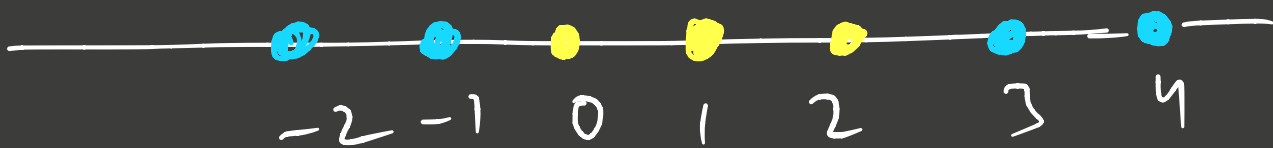
$$K(x_i, x_j) = (1 + x_i x_j)^2 - 1$$

↑ dot product in lower dimension

$$= 1 + 2x_i x_j + x_i^2 x_j^2 - 1$$

$$= \langle \sqrt{2}x_i, x_i^2 \rangle \cdot \langle \sqrt{2}x_j, x_j^2 \rangle$$

$$= \phi(x_i) \cdot \phi(x_j)$$



## Original Dataset

#	$x$	$y$
1	-2	-1
2	-1	-1
3	0	1
⋮	⋮	⋮

## Transformed dataset

#	$\sqrt{2}x$	$x^2$	$y$
1	$-2\sqrt{2}$	4	-1
2	$-\sqrt{2}$	1	-1
3	0	0	1
⋮	⋮	⋮	⋮

$$\phi(x_1) = \langle -2\sqrt{2}, 4 \rangle; \quad \phi(x_2) = \langle -\sqrt{2}, 1 \rangle \quad \text{TRANSFORMAT}^N$$

$$\phi(x_1) \cdot \phi(x_2) = \underline{-2\sqrt{2}} * \underline{-\sqrt{2}} + \underline{4} * \underline{1} = \underline{8} \quad \text{DOT PRODUCT IN 2D}$$

$$k(x_1, x_2) = \left\{ 1 + \underline{(-2) * (-1)} \right\}^{-1} \quad \text{DOT PRODUCT IN 1D}$$



WHY DID WE USE DUAL FORM?

KERNELS AGAIN!!

PRIMAL FORM DOESN'T ALLOW

FOR "KERNEL TRICK"

$K(\vec{x}_1, \vec{x}_2)$  in DUAL

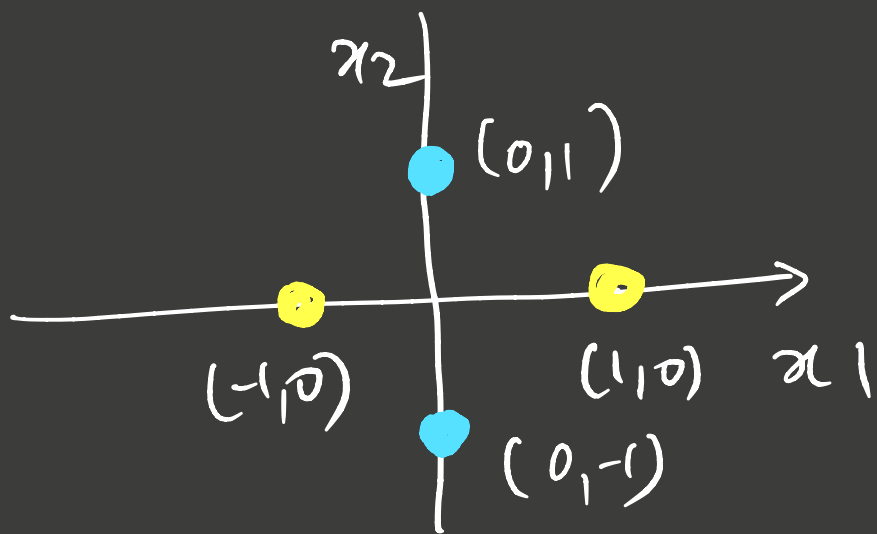
↳ COMPUTE  $\phi(x)$  and then dot product in 'D' dimensions?

GRAM MATRIX (Positive Semi-Definite)

$$K(x_i, x_j) = (1 + x_i \cdot x_j)^2$$

	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$
$x_1$	24	8	0	0	8	24	48
$x_2$	8	1	0	-1	0		
$x_3$	0	..	..	..	..		
$x_4$	0						
$x_5$	8						
$x_6$	24						
$x_7$	48						

# ANOTHER EXAMPLE



$$K(\bar{x}, \bar{x}') = (\vec{x}^T \cdot \vec{x}')^2$$

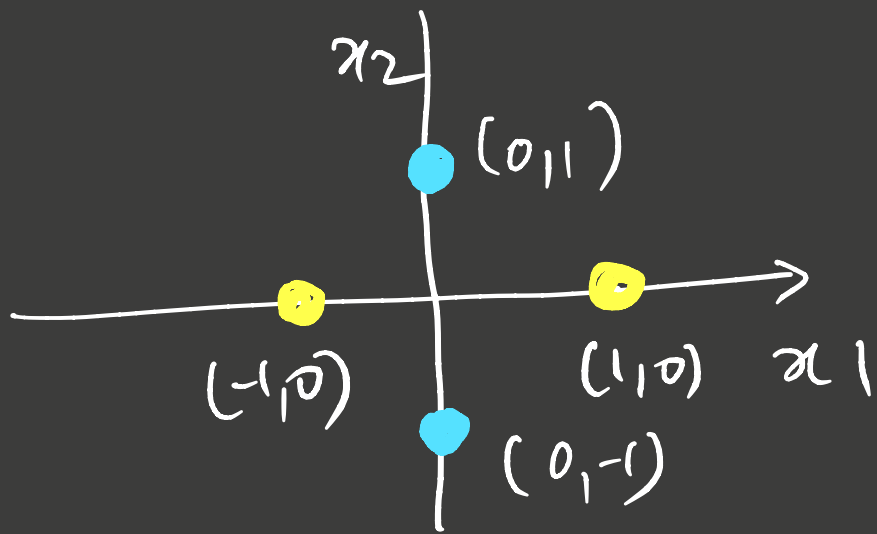
$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}; \bar{x}' = \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix}$$

Q: What is  $\phi(x)$ ?

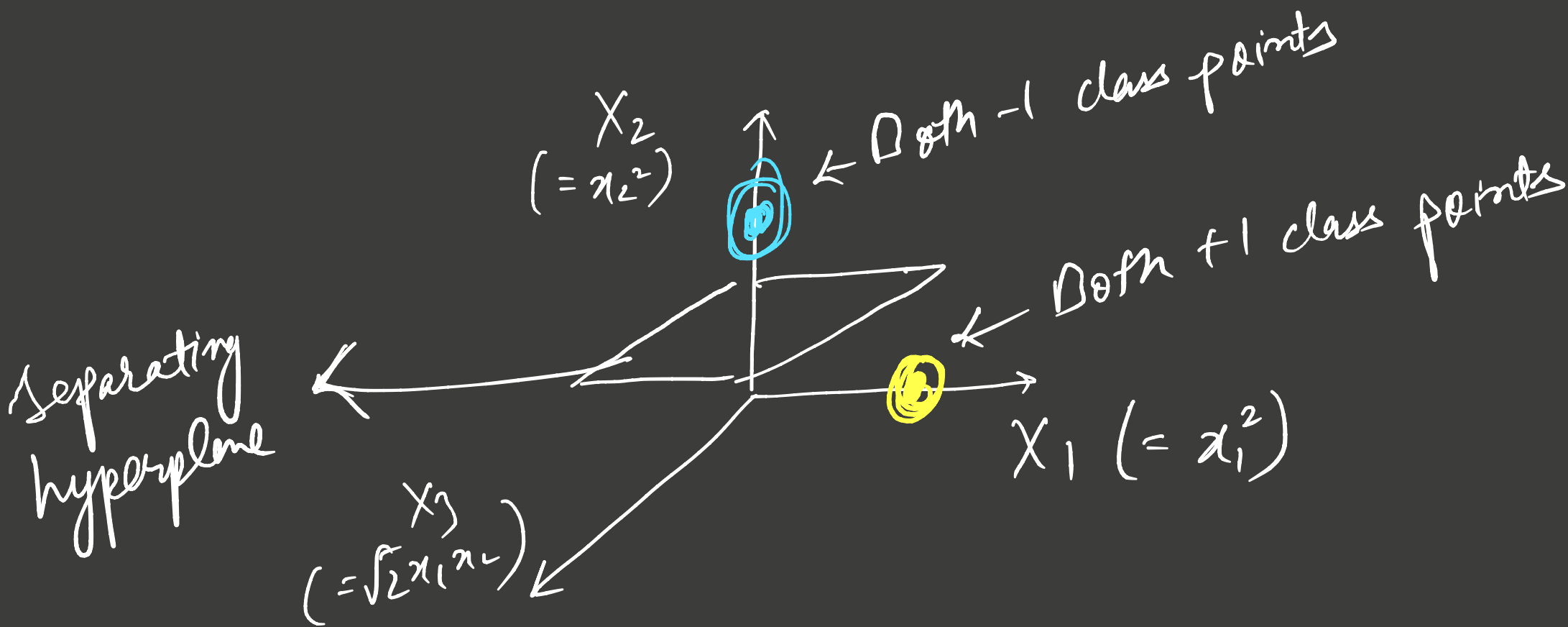
$$K(\bar{x}, \bar{x}') = \phi(\bar{x}) \cdot \phi(\bar{x}')$$

$$K(\bar{x}, \bar{x}') = \left\{ \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} x'_1 \\ x'_2 \end{bmatrix} \right\}^2 = (x_1 x'_1 + x_2 x'_2)^2$$

$$\Rightarrow \boxed{\phi(x) = \langle x_1^2, \sqrt{2}x_1x_2, x_2^2 \rangle} = x_1^2 x_1' + x_2^2 x_2' + 2x_1 x_1' x_2 x_2'$$



$\Downarrow \phi(x)$



# SOME KERNELS

① Linear:  $K(\bar{x}_1, \bar{x}_2) = \bar{x}_1 \cdot \bar{x}_2$

② Polynomial:  $K(\bar{x}_1, \bar{x}_2) = (p + \bar{x}_1 \cdot \bar{x}_2)^q$

③ Gaussian:  $K(\bar{x}_1, \bar{x}_2) = e^{-\gamma \|\bar{x}_1 - \bar{x}_2\|^2}$

ALSO CALLED RADIAL BASIS FUNCTION (RBF)

$$\gamma = \frac{1}{2\sigma^2}$$

0) For  $\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$  what space does

Kernel  $K(\bar{x}, \bar{x}') = (1 + \bar{x} \cdot \bar{x}')^3$  belong to?

or

$$\bar{x} \in \mathbb{R}^2$$

$$\phi(\bar{x}) \in \mathbb{R}^?$$

$$K(x, z) = (1 + x_1 z_1 + x_2 z_2)^3$$

$$= \dots = \langle 1, x_1, x_2, x_1^2, x_2^2, x_1^2 x_2, x_1 x_2^2, x_1^3, x_2^3 \rangle$$

10 dimensional!

0) For  $\bar{x} = x$ ; what space does RBF kernel lie in?

$$K(x, z) = e^{-\gamma \|x - z\|^2}$$
$$= e^{-\gamma (x - z)^2}$$

Now;  $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$

$\therefore e^{-\gamma (x - z)^2}$  is  $\infty$  dimensional !!

Q)  $I_1$  SUM parametric or non-parametric?



Q) Is SVM parametric or non-parametric?

Yes and No  
↓ ↓

Linear

kernel

Polynomial  
kernel

(form fixed)

RBF

(form changes with  
data)

RBF is Non-Parametric

$$\hat{y}(x_{\text{Test}}) = \text{SIGN}(\vec{w} \cdot \vec{x}_{\text{Test}} + b)$$

$$= \text{SIGN}\left(\sum_{j=1}^{N_{\text{sv}}} \alpha_j y_j \vec{x}_j \cdot \vec{x}_{\text{Test}} + b\right)$$

⇓ Kernelized.

$$\hat{y}(x_{\text{Test}}) = \text{SIGN}\left(\sum_{j=1}^N \alpha_j y_j K(\vec{x}_j, \vec{x}_{\text{Test}}) + b\right)$$

$\alpha_j = 0$  where  $j \neq \text{S.V.}$

New  $K(\vec{x}_j, \vec{x}_{\text{Test}})$  for RBF is:

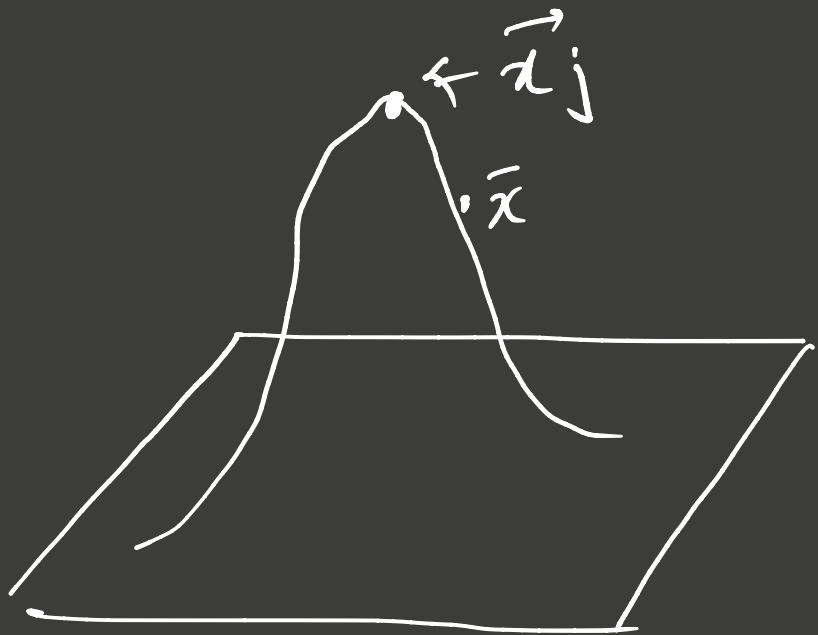
$$\frac{-\gamma \|\vec{x}_j - \vec{x}_{\text{Test}}\|^2}{e}$$

$\therefore$  Hypothesis is a function of "All" train points.

↓  
What kind of?

Close  $\vec{x}$  is to  $\vec{x}_j$ , more is it influencing  $\hat{y}(\vec{x})$

← Hypothesis function



Now if we add a point to  
dataset



Functional form can  
adapt (similar to  
KNN)

∴ SUM with RBF Kernel

$\hat{f}_\Delta$  Non-Parametric

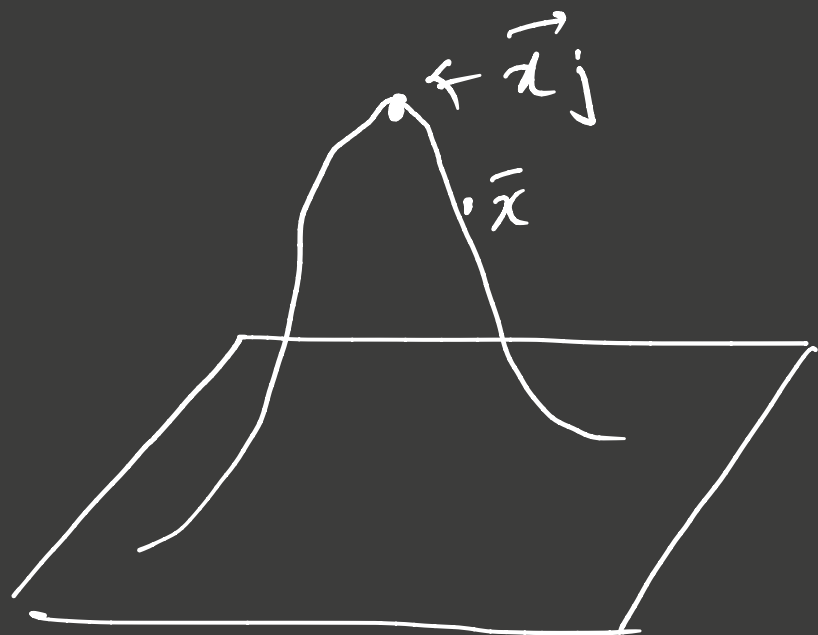
# Interpretation of RBF

$$\hat{y}(x) = \text{SIGN} \left( \underbrace{\sum \alpha_i y_i}_{\text{Activation}} \underbrace{e^{-\|x - x_i\|^2}}_{\text{Basis}} + b \right)$$

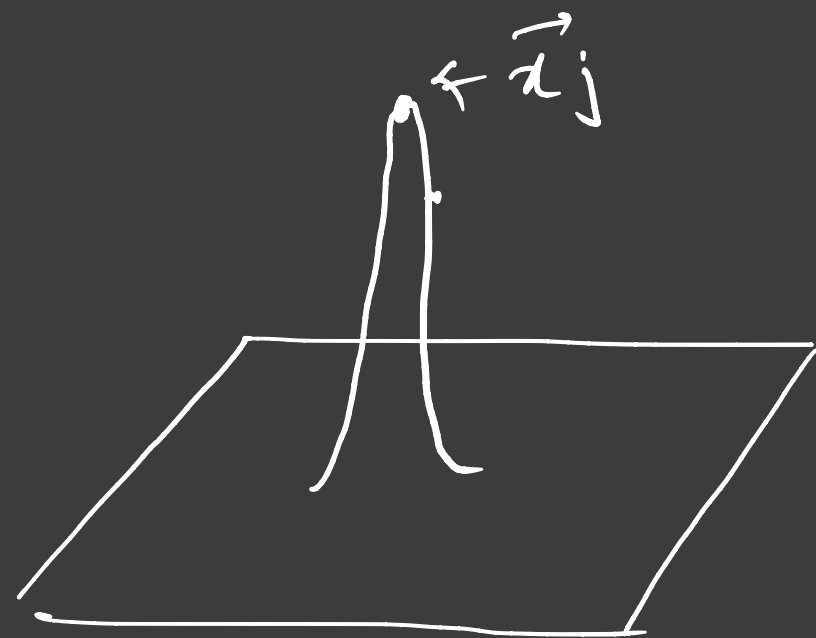
↙ Radial

# RBF : Effect of $\gamma$

$\gamma$ : How far is the influence of a single training sample



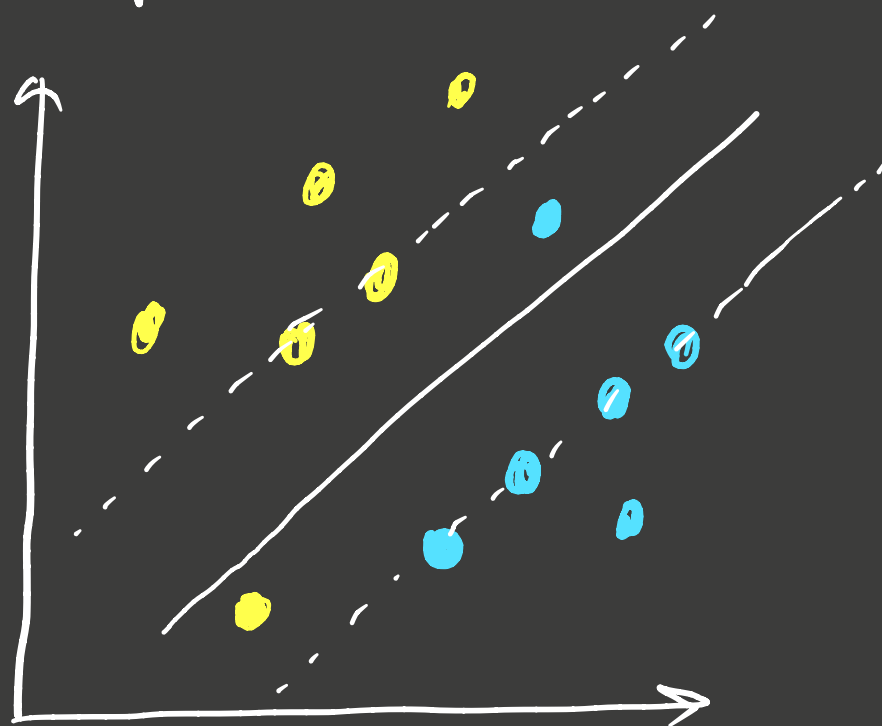
$\gamma = \text{Low}$   
High influence of  $\bar{x}_j$



$\gamma = \text{High}$   
low influence of  $\bar{x}_j$

# SOFT MARGIN SVM

Q: Can we learn SVM for "slightly" non separable data without projecting to a higher space?

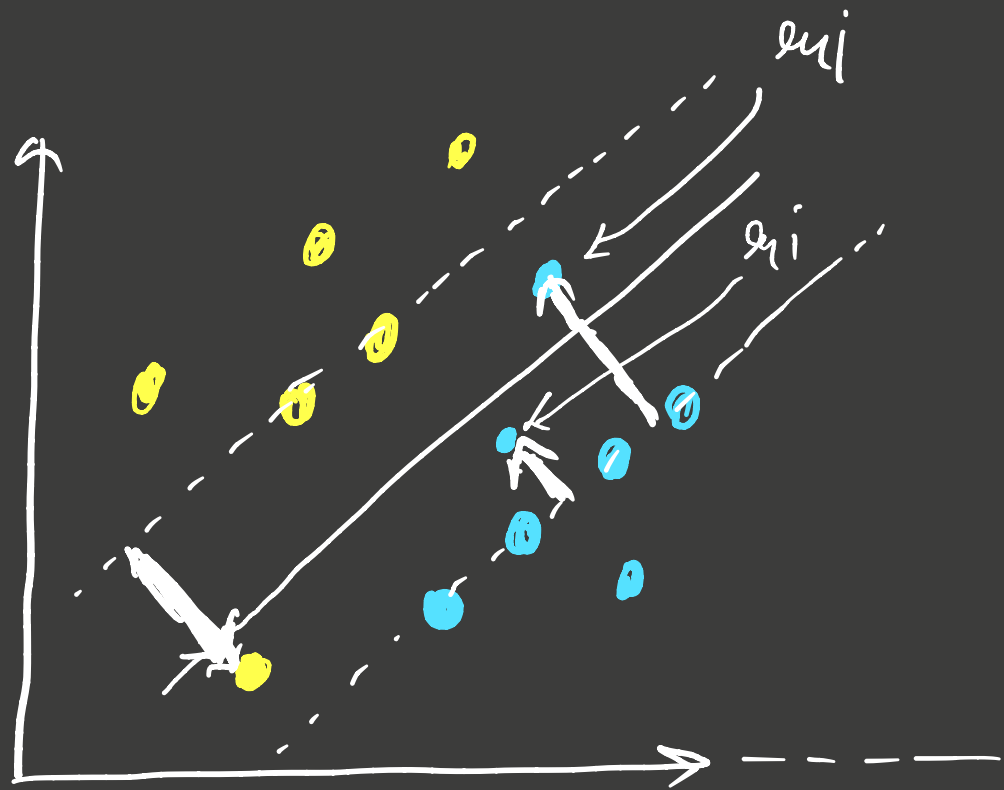


"Slightly"  
non  
separable  
data.

# SOFT MARGIN SVM

SLACK VARIABLE

$$\xi_i = \begin{cases} 0 & \text{if point on correct side of margin} \\ \text{Distance from margin} & \text{otherwise} \end{cases}$$



Change objective

$$\min \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 - \xi_i; \xi_i \geq 0;$$

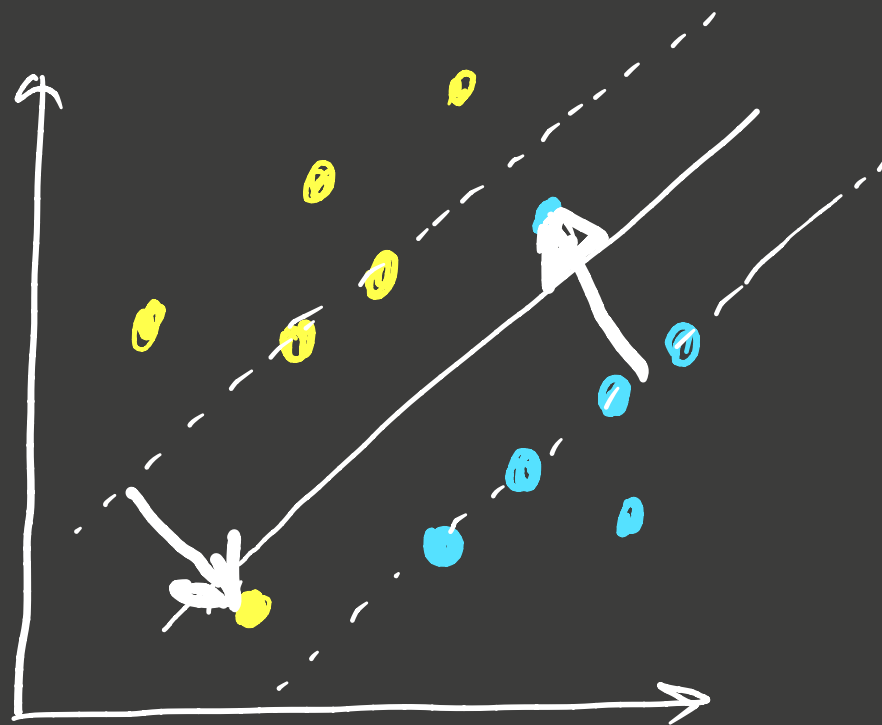


# SOFT MARGIN SVM

Change objective

$$\min \frac{1}{2} \|\bar{w}\|^2 + C \sum_{i=1}^n \eta_i$$

$$\text{s.t. } y_i (\bar{w} \cdot \bar{x}_i + b) \geq 1 - \eta_i$$



in Dual

$$\text{Maximize } \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \bar{x}_i \cdot \bar{x}_j$$

s.t.

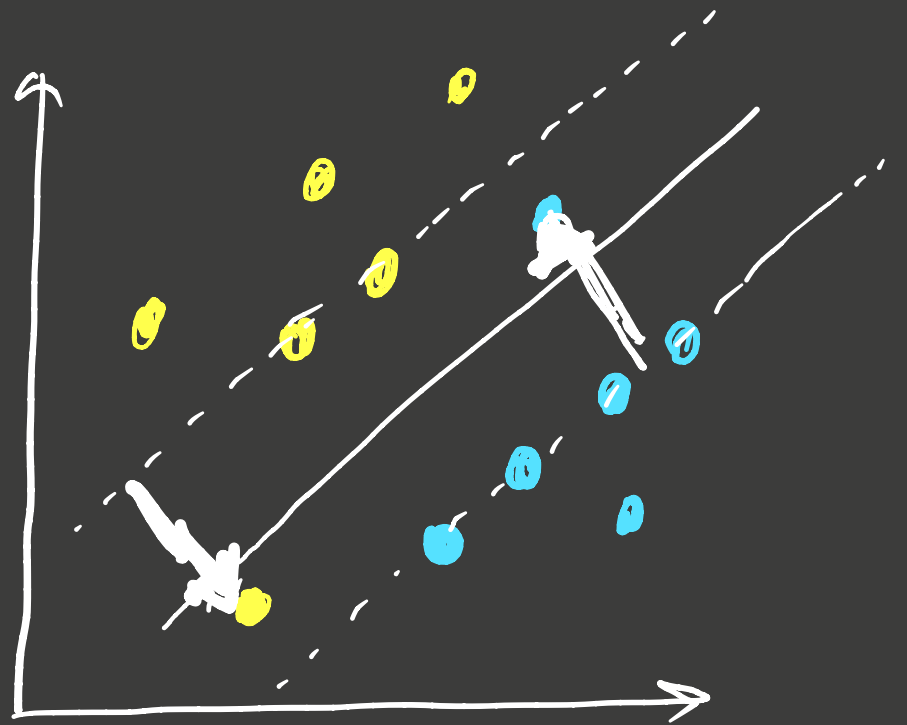
$$0 \leq \alpha_i \leq C \quad \& \quad \sum_{i=1}^n \alpha_i y_i = 0$$

# SOFT MARGIN SVM

Change objective

$$\min \frac{1}{2} \|\bar{w}\|^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y_i (\bar{w} \cdot \bar{x}_i + b) \geq 1 - \xi_i$$

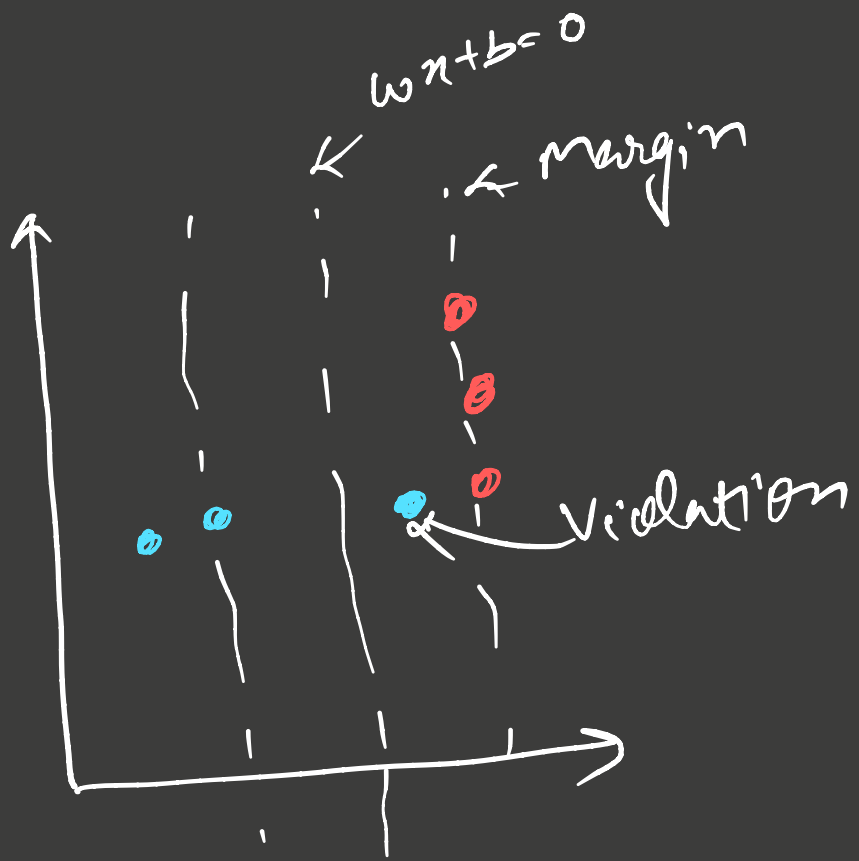


$C \rightarrow 0$  : Larger margin

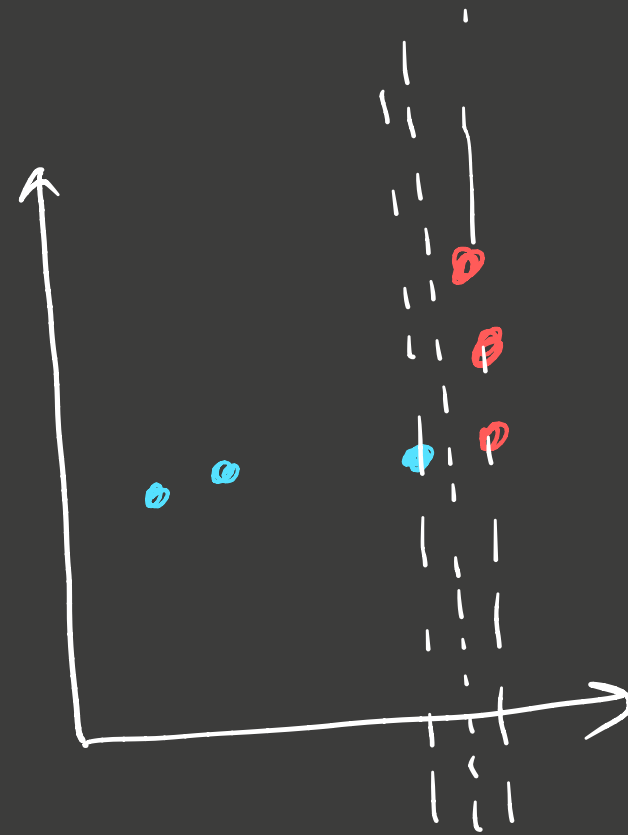
$C \rightarrow \infty$  : Smaller margin

why?

Notes ok : SVM - soft-margin



Low value of  $C$



High value of  $C$

If  $C \rightarrow 0$

Objective  $\rightarrow \min \frac{1}{2} \|\vec{w}\|^2$

$\Rightarrow$  Choose large margin

(without worrying for  $\epsilon_i$ 's)

[ Recall : Margin =  $\frac{2}{\|\vec{w}\|}$  ]

If  $C \rightarrow \infty$  (or very large)

Objective  $\rightarrow \min C \sum \epsilon_i^2$  or choose  
' $w, b$ ', s.t.  $\epsilon_i$  is small.

Q) what is equivalent of hard margin?

a)  $C \rightarrow 0$

b)  $C \rightarrow \infty$

Q) what is equivalent of hard margin?

$$a) C \rightarrow 0$$

$$b) C \rightarrow \infty$$



No violations!!

BIAS      VARIANCE      TRADE-OFF      FOR  
SOFT-MARGIN  
SVM

Low  $C \Rightarrow$  Higher train error  
(higher bias)

High  $C \Rightarrow$  Very sensitive to dataset  
(high variance)

# SOFT MARGIN SVM

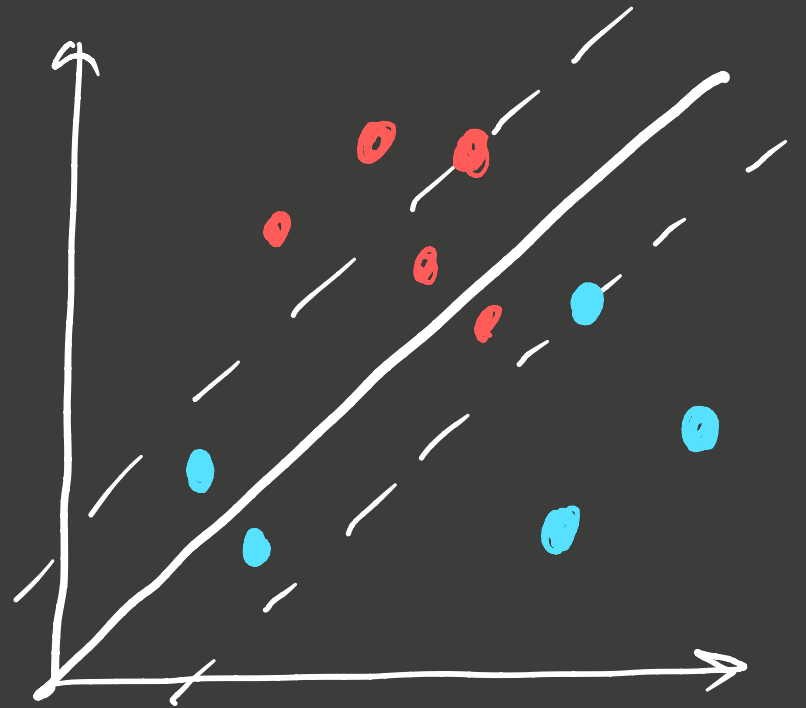
Types of support vectors

Zone 1  $y_i (\bar{w} \cdot \bar{x}_i + b) = 1$

Zone 2  $0 < y_i < 1$  (correctly classified)

Zone 3  $y_i > 1$  (misclassified)

$\therefore$  As  $C$  increases, # support vectors decreases



Notes ok: SVM - soft-margin



# SVM FORMULATION IN LOSS + PENALTY FORM

---

Objective:

$$\min \frac{1}{2} \|\bar{w}\|^2 + C \sum_{i=1}^N \xi_i$$

$$\text{s.t.} \quad y_i (\bar{w} \cdot \bar{x}_i + b) \geq 1 - \xi_i \quad ; \quad \xi_i \geq 0$$

$$\text{New:} \quad y_i (\bar{w} \cdot \bar{x}_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 1 - y_i (\bar{w} \cdot \bar{x}_i + b)$$

But  $\xi_i \geq 0$

$$\therefore \xi_i = \text{MAX} [0, 1 - y_i (\bar{w} \cdot \bar{x}_i + b)]$$

∴ Objective is:

$$\text{Min } C \sum y_i + \frac{1}{2} \|\vec{w}\|^2$$

$$\text{or Min } C \sum_{i=1}^N \text{MAX} [0, 1 - y_i (\vec{w} \cdot \vec{x}_i + b)] + \frac{1}{2} \|\vec{w}\|^2$$

$$\text{or Min } \sum_{i=1}^N \text{MAX} [0, 1 - y_i (\vec{w} \cdot \vec{x}_i + b)] + \frac{1}{2C} \|\vec{w}\|^2$$

Loss

REGULARISATION

# LOSS FUNCTION FOR SUM (HINGE LOSS)

LOSS FUNCTION IS:  $\sum_{i=1}^N \text{MAX}[0, 1 - y_i (\bar{w} \cdot \bar{x}_i + b)]$

Case I

$$y_i (\bar{w} \cdot \bar{x}_i + b) = 1$$

LIES ON MARGIN LOSS  $i$  = 0

Case II

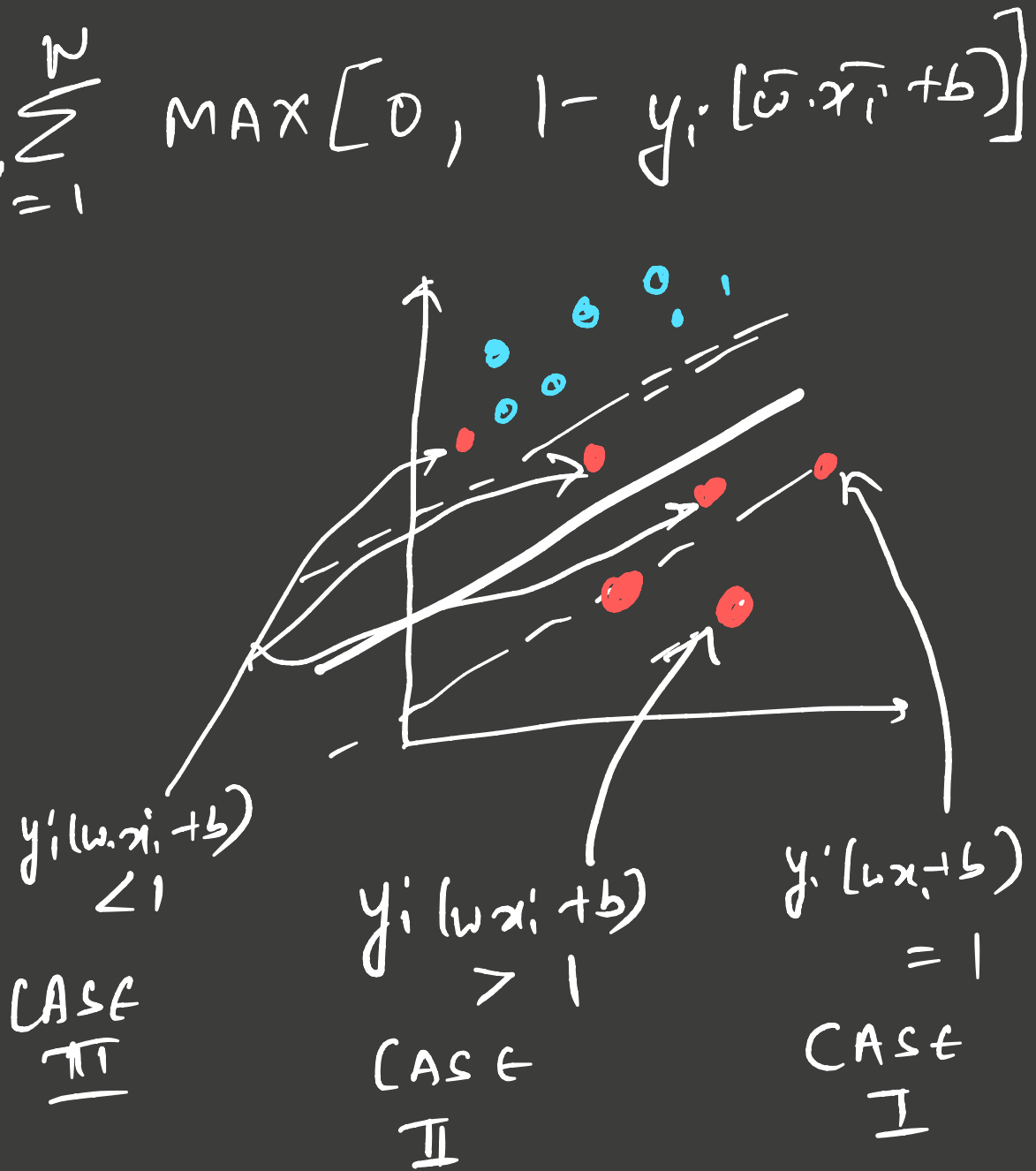
$$y_i (\bar{w} \cdot \bar{x}_i + b) > 1$$

$$\text{Loss}_i = 0$$

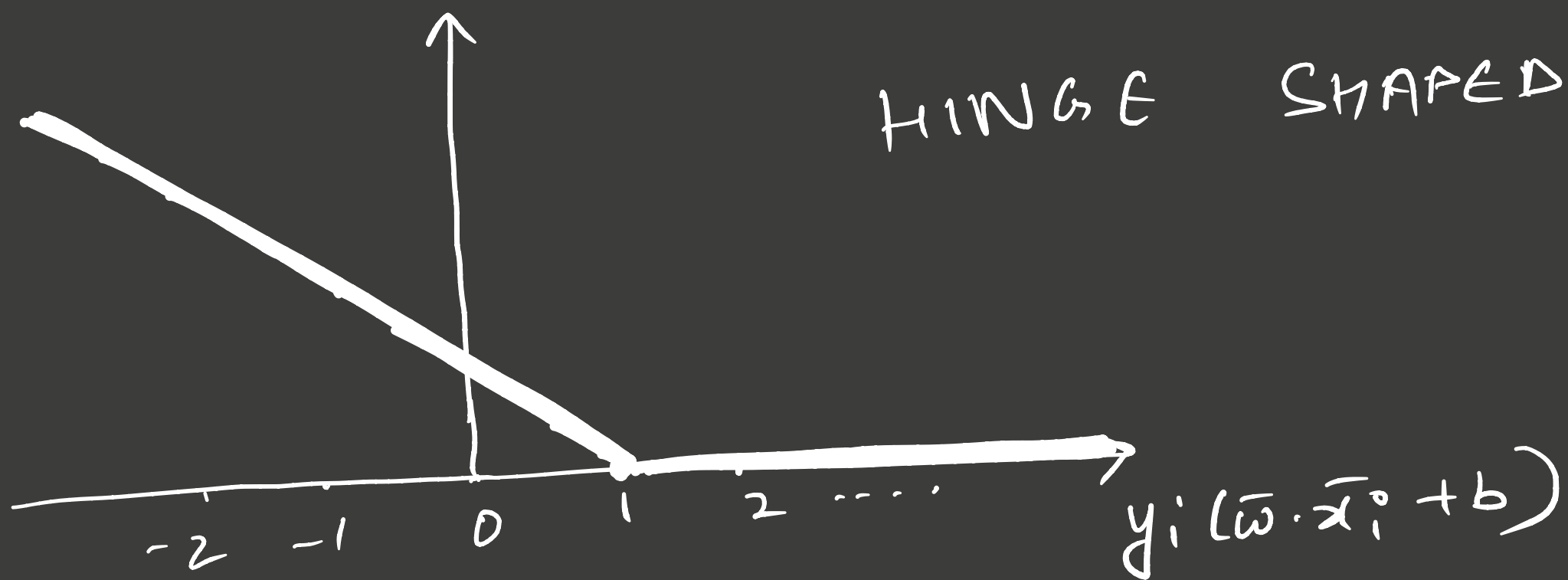
Case III

$$y_i (\bar{w} \cdot \bar{x}_i + b) < 1$$

$$\text{Loss}_i \neq 0$$



# HINGE LOSS CONTINUED



Q) IS HINGE LOSS CONVEX & DIFFERENTIABLE?

CONVEX: ✓

DIFFERENTIABLE: ✗

SUBGRADIENT: ✓

SUM LOSS IS CONVEX

HINGE LOSS  $\sum (\max [0, (1 - y_i (\bar{w} \cdot x_i + b))])$   
IS CONVEX

PENALTY  $\frac{1}{2} \|\bar{w}\|^2$   
IS CONVEX

$\therefore$  SUM LOSS IS CONVEX