

24/9/18

- we were interested in a hyper plane that separates the classes.

- The plane which places training data in the correct location

- we want the hyper plane, that also maximizes distance b/w the classes

- A normal classifier gives a feasible solution, but not enough in practice

- SVM can do that  $\rightarrow$  (that  $\Rightarrow$  it can classify & ensure maximum distance b/w points)

- ~~Draw 2 lines,~~

- Now our problem is maximize margin such that samples are correctly classified

- $\pm 1$  don't have a significance.

Distance b/w the separators

$$\frac{b+1}{\|w\|} - \frac{b-1}{\|w\|} = \frac{2}{\|w\|}$$

needs to be maximized

and  $\begin{cases} w^T x_i + b \geq 1 & \text{if } y_i = +1 \\ w^T x_i + b \leq -1 & \text{if } y_i = -1 \end{cases}$

Combine

$$y_i(w^T x_i + b) \geq 1, \forall i$$

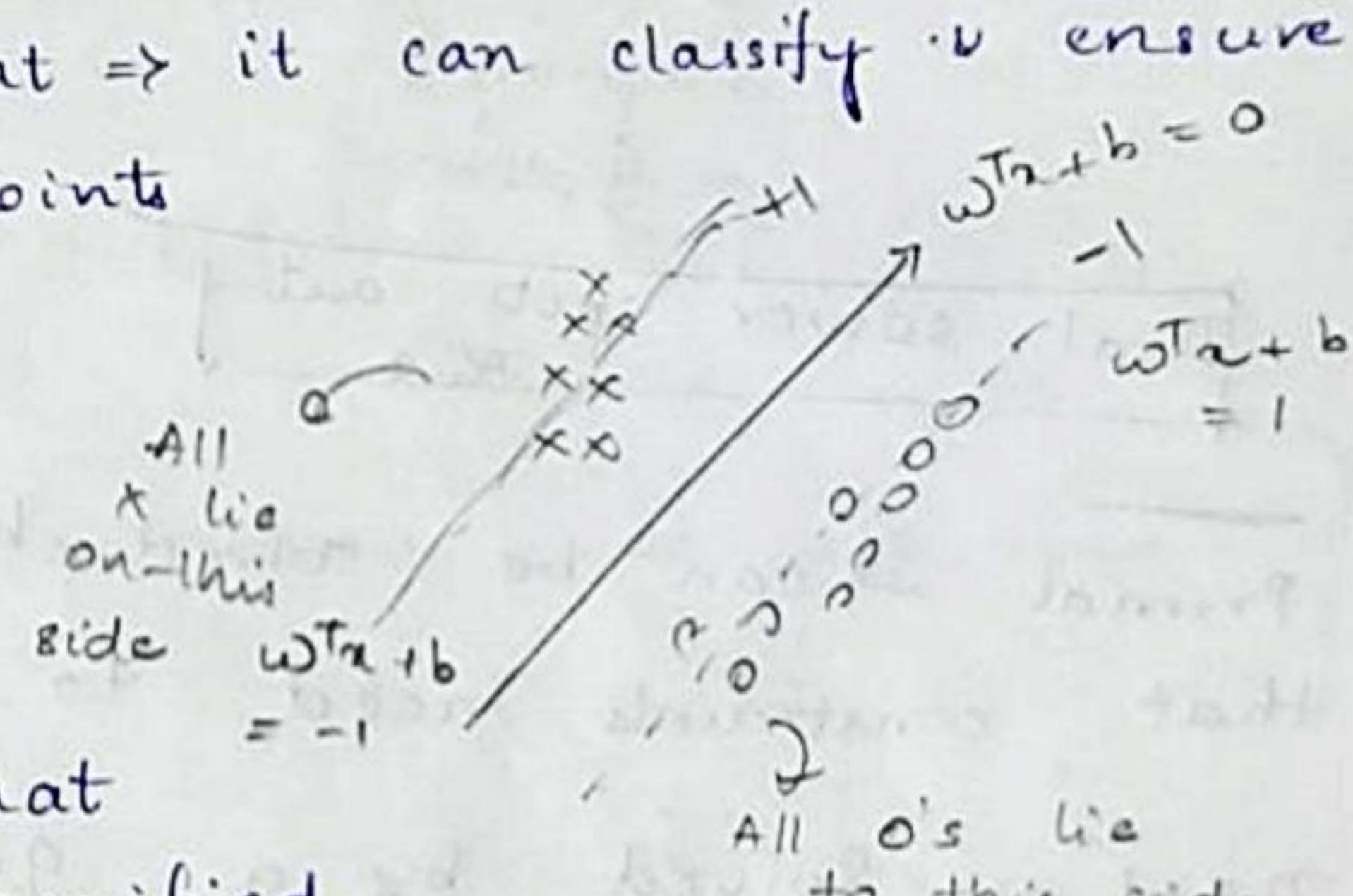
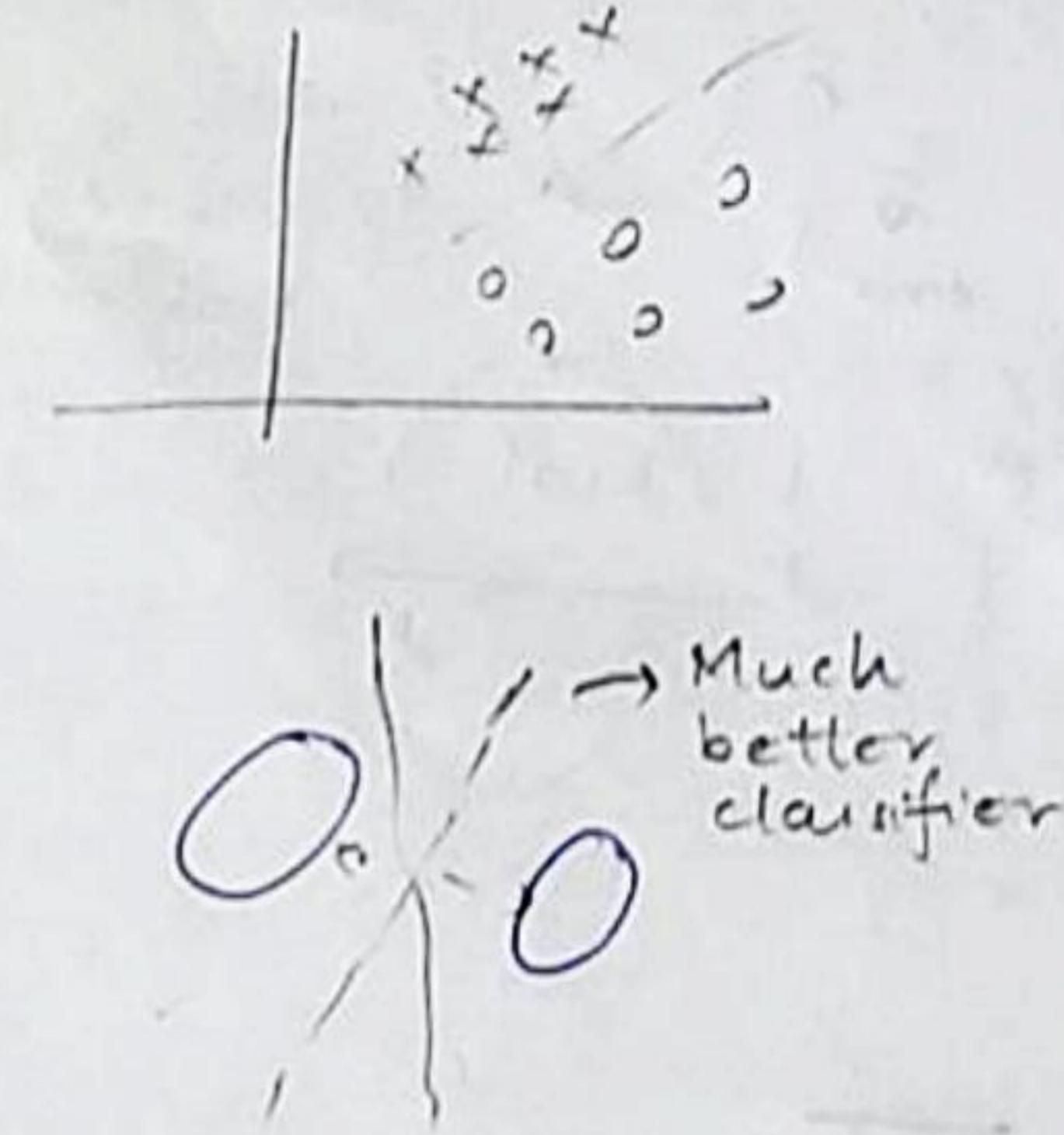
Converting

minimize  $\rightarrow \frac{\|w\|}{2} = \frac{w^T w}{2}$ ; with condition

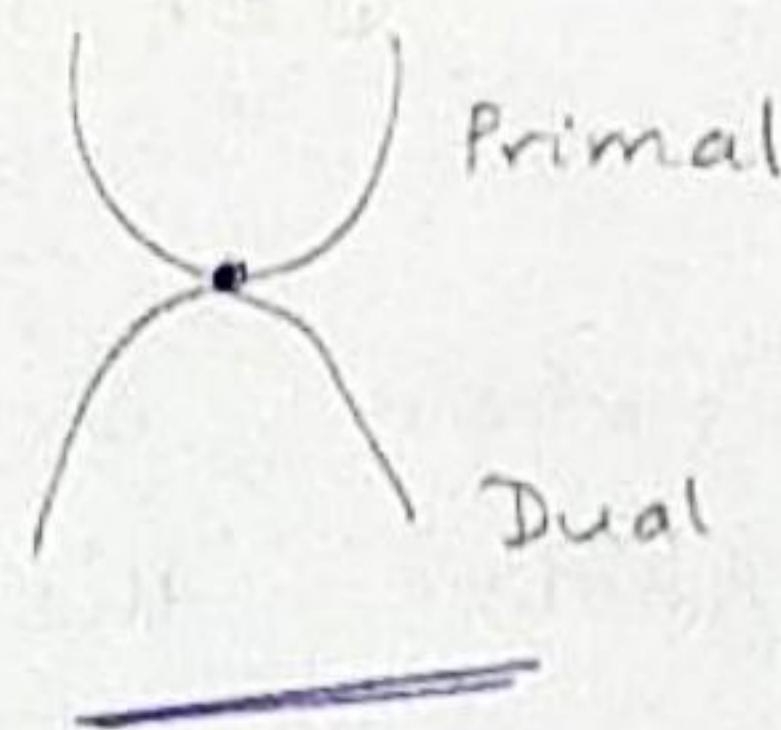
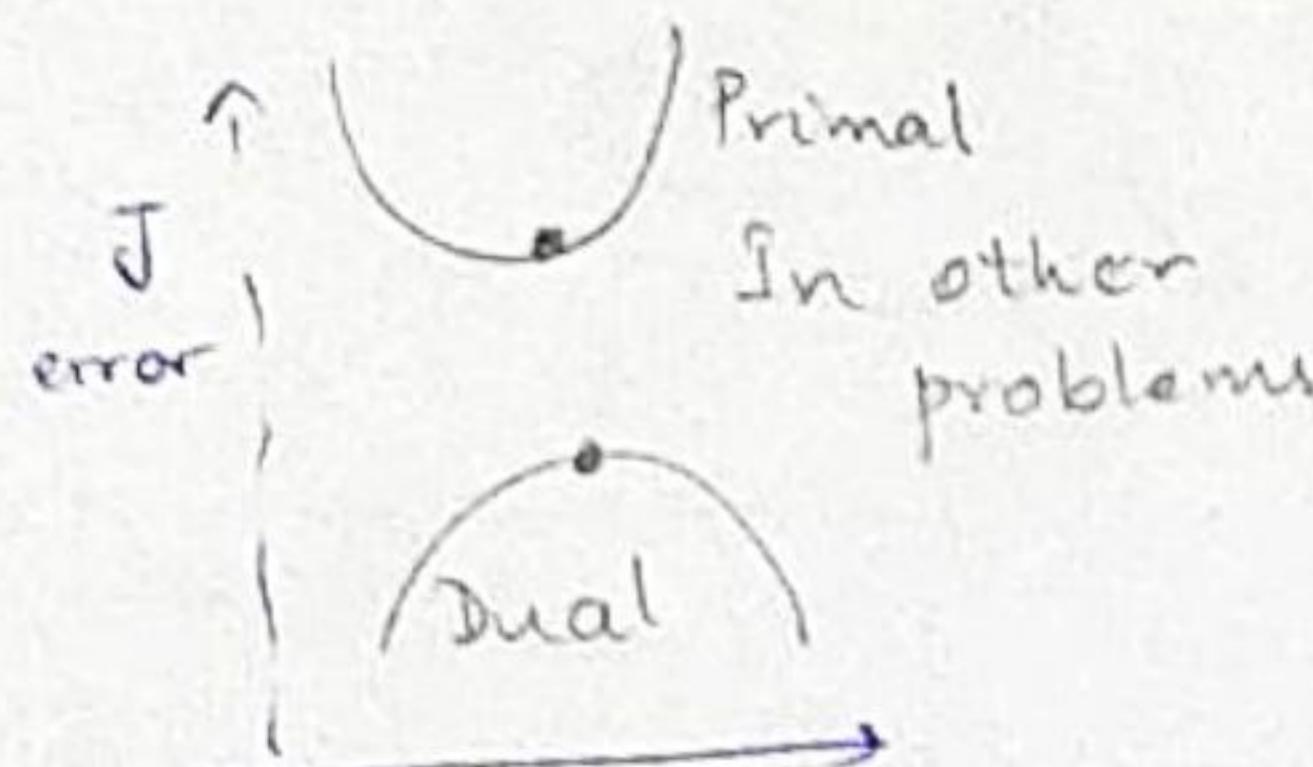
$\downarrow$  primal

$$\begin{aligned} y_i(w^T x_i + b) &\geq 1 \\ y_i \in \{+1, -1\} \end{aligned}$$

- SVM is actually practised as per the theoretical basis //



## Dual x Primal



Both have the same error, as SVM is a well-behaving problem.

Dual: Max  $\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (\alpha_i \alpha_j y_i y_j x_i^T x_j)$

Lagrangian  
multiplier

Dual solver spits out  $\alpha_i$ 's

$$\sum_{i=1}^N \alpha_i y_i = 0 \quad \alpha_i \geq 0$$

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

Primal → can be solved by Q.D., (keeping in mind that constraints need to be satisfied).

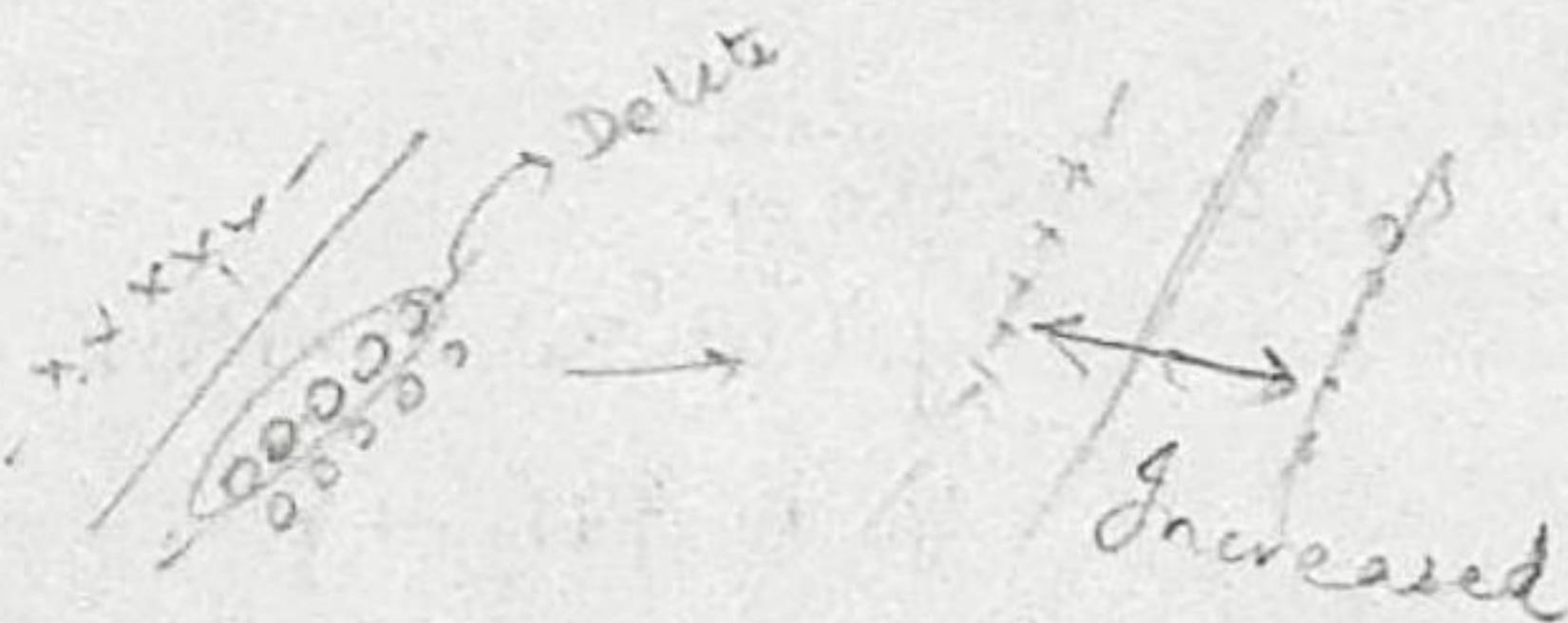
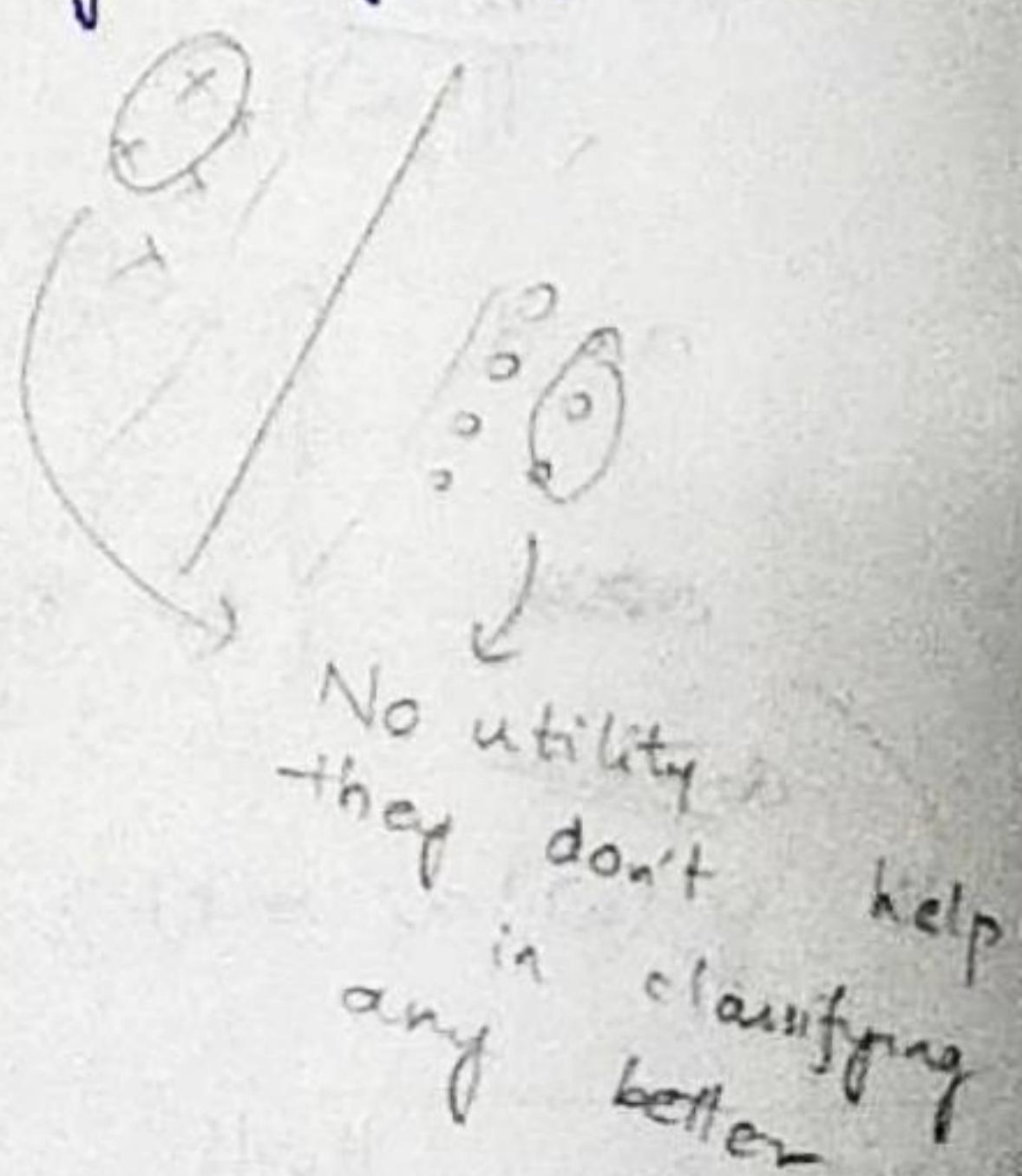
Dual → Solved by a quadratic solver

• Why dual is needed -? : We can appreciate the soln better.

• SVM looks for Max-margin separating hyperplane

•  $\alpha = 0$ ; for samples that don't matter.

So  $\alpha$  is non zero for the frontier points



So  $\alpha$  is sparse  $\alpha = [\alpha_i \text{ from } 1 \text{ to } n] = [\alpha_1, \alpha_2, \dots, \alpha_n]$   
i.e. for most samples,  $\alpha$  is ~~sparse or zero~~.

— Most solvers give non-zero  $\alpha$ 's.  
If  $\text{#sv} \cdot \alpha < n \Rightarrow$

support vectors  $\Rightarrow$  those points with non-zero  $\alpha$ .

The final decisions depends on sv?

$$w = \sum_{i=1}^N \alpha_i y_i x_i \xrightarrow{\text{convert}} w = \sum_{i \in \text{sv}} \alpha_i y_i x_i$$

Toy prob

- If 1M samples & 100 SVMs  $\Rightarrow$  Why do you need others?  
don't
- ~~first~~ We ~~need~~ know SVM's in the first place
- Only after all shit is done, we get soln.

Testing data:

$$\begin{aligned} w^T x + b &\leq 0 \\ \sum \alpha_i y_i x_i^T x + b &\leq 0 \end{aligned} \quad \text{Put } w.$$

Notice that everywhere  $x^T x$  is present. ~~so~~,  
we will exploit it, later

Leave-1-out

Remove 1 sample & build the classifier  
Test on it & find avg

In SVMs :-  
If non sv  $\Rightarrow$  no error

If sv  $\Rightarrow$  then, maybe error.

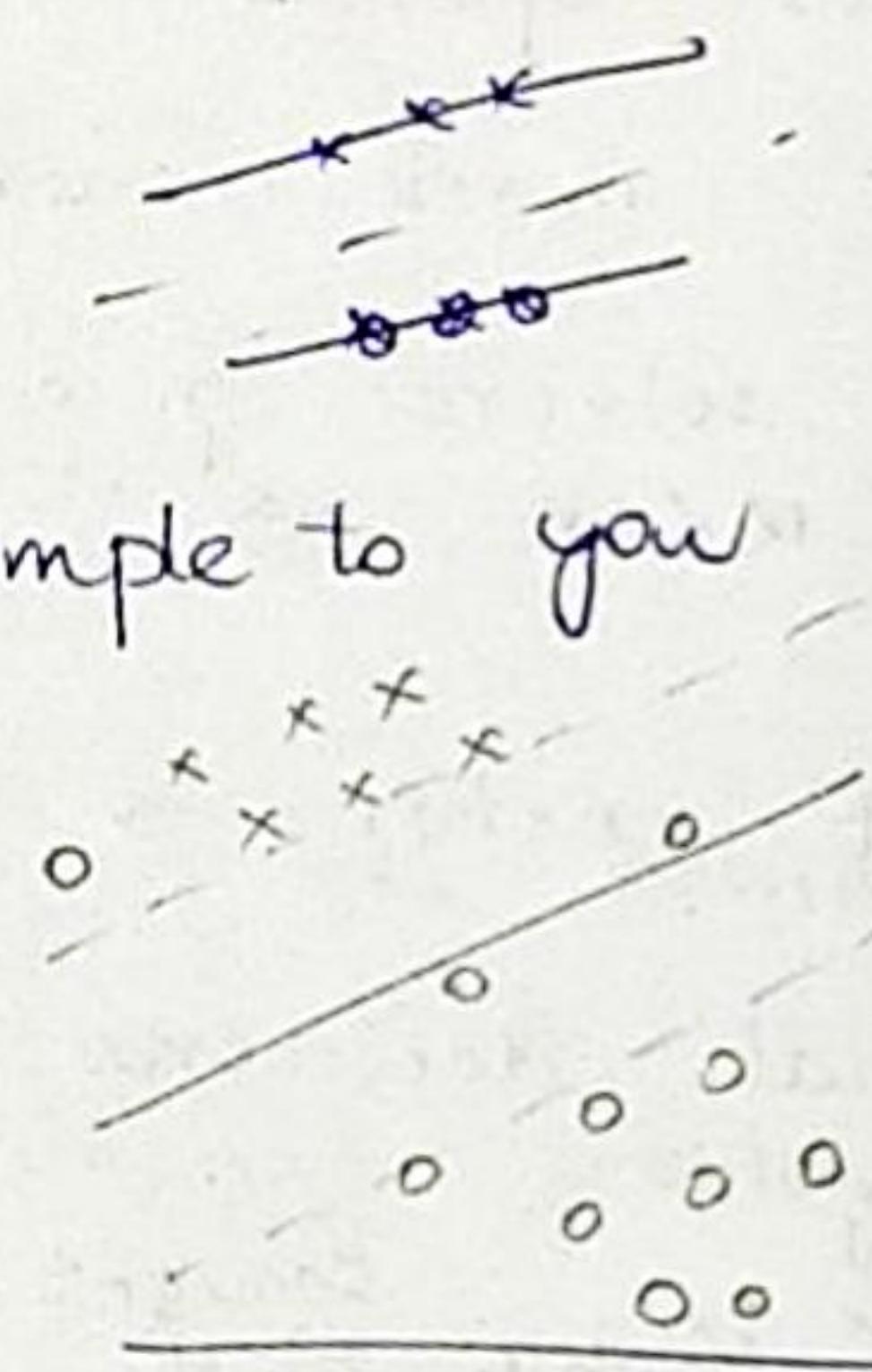
So,  $\frac{\# \text{sv}}{N}$  is an upperbound on the error

- If  $\omega$ ,  $b$  are present,  
why do you need other SVs?

Ans: It is in 2D, so it looks simple to you  
 $\downarrow$

Explanation to this stupid ans:

- If suppose the classification looks like this: (Samples can never be linearly separable)  
 We want  $\epsilon_i$  to be minimized] (separable with a line)



$$\left[ \begin{array}{l} \omega^T x_i + b \geq +1 - \epsilon_i \text{ for } y_i = +1 \\ \epsilon_i \geq 0 \end{array} \right]$$

So minimize  $\frac{1}{2} \omega^T \omega + \sum_{i=1}^N \epsilon_i$

- Prev. prob: Hard margin SVM,
- This. prob: Soft margin SVM

Append 1  $\frac{1}{2} \omega^T \omega + C \sum_{i=1}^N \epsilon_i$

If  $C$  small  $\Rightarrow$  anyone can violate  
 $\epsilon_i$  needs to be optimized

Append 2  $\frac{1}{2} \omega^T \omega + C \times \sum_{i=1}^N \epsilon_i^k \rightarrow L_k - \text{SVM}$

- Minimum number of penalty
- Minimum amount of penalty



for some  $k$

$\hookrightarrow$  strictly convex  $\rightarrow$

for others  $\hookrightarrow$  convex  $\rightarrow$

Back to primal

$$\frac{1}{2} \omega^T \omega - \sum_{i=1}^N \alpha_i [y_i (\omega^T x_i)]$$

$$J(\omega, b, \alpha) = \frac{1}{2} \omega^T \omega - \sum_{i=1}^N \alpha_i [y_i (\omega^T x_i + b) - 1]$$

2 things  $\Rightarrow \alpha_i = 0$  or // not on wall  
 $y_i (\omega^T x_i + b) = 1 \Rightarrow // \text{on wall}$

$$\frac{\partial J}{\partial \omega} = 0, \quad \frac{\partial J}{\partial b} = 0, \quad \rightarrow \text{Given}, \quad \sum_{i=1}^N \alpha_i y_i = 0.$$

$$\hookrightarrow \text{Given } \omega = \sum_{i=1}^N \alpha_i y_i \omega^T x_i$$

$$\begin{aligned} J &= \frac{1}{2} \left[ \sum_{i=1}^N \alpha_i y_i x_i^T \right] \left[ \sum_{j=1}^N \alpha_j y_j x_j^T \right] - \sum_{i=1}^N \alpha_i [y_i \omega^T x_i] \\ &\quad - \sum_{i=1}^N \alpha_i y_i b + \underbrace{\sum_{i=1}^N \alpha_i}_{0} y_i = 0 \end{aligned}$$

$$= \frac{1}{2} \left[ \sum_{i=1}^N \alpha_i y_i x_i^T \right] \left[ \sum_{j=1}^N \alpha_j y_j x_j^T \right] - \sum_{i=1}^N \alpha_i \underbrace{[y_i \omega^T x_i]}_{\downarrow \text{Substitute } \omega} + \sum_{i=1}^N \alpha_i$$

$$= \frac{1}{2} \left[ \sum_{i=1}^N \alpha_i y_i x_i^T \right] \left[ \sum_{j=1}^N \alpha_j y_j x_j^T \right] - \left[ \sum_{i=1}^N [\alpha_i y_i \omega^T x_i] \right] \left[ \sum_{i=1}^N \alpha_i y_i x_i^T \right] + \sum_{i=1}^N \alpha_i$$
$$\frac{1}{2} \cancel{\alpha} - \cancel{\alpha} = -\frac{1}{2} \cancel{\alpha}$$

$$= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i^T x_j^T + \sum_{i=1}^N \alpha_i$$
$$\cancel{\alpha}$$

27/9/18

Lec-15

## Kernels and Non-linear SVM's

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

If  $d \leq N$

$$\text{and } \sum \sum \alpha_i \alpha_j y_i y_j x_i^T x_j = \sum_{i=1}^N \alpha_i w^T x_i \rightarrow O(d)$$

If  $w^T x$  is already known  
computation

Else  $\sum \alpha_i \alpha_j y_i y_j x_i^T x_j \rightarrow d \times \# \text{SV}$

$$x \Rightarrow x' = Mx \quad \text{Linear dimensionality reduction}$$

- Kernels do a transformation

$x \rightarrow \underline{\phi(x)}$  generally, it increases dimensionality

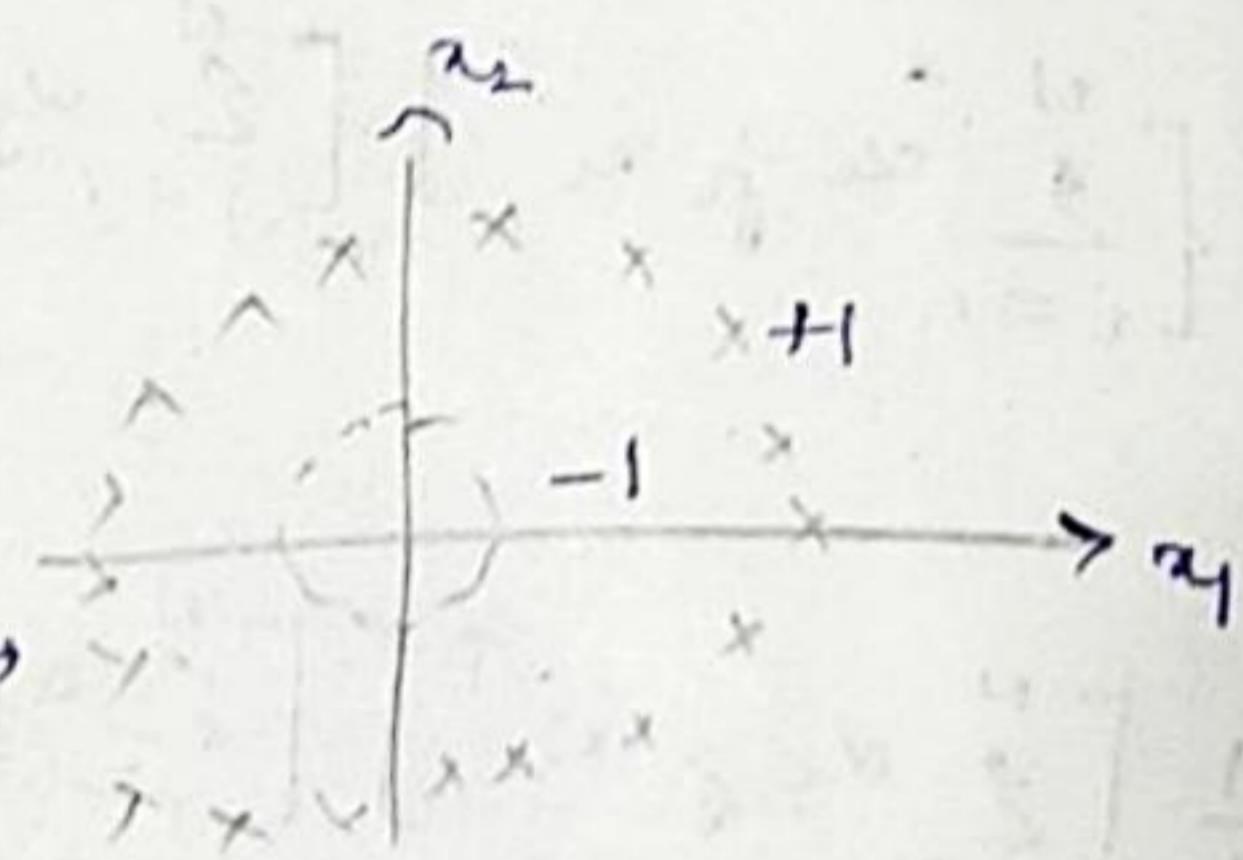
↓  
Some non-linear mapping

A linear hyperplane is not separating.

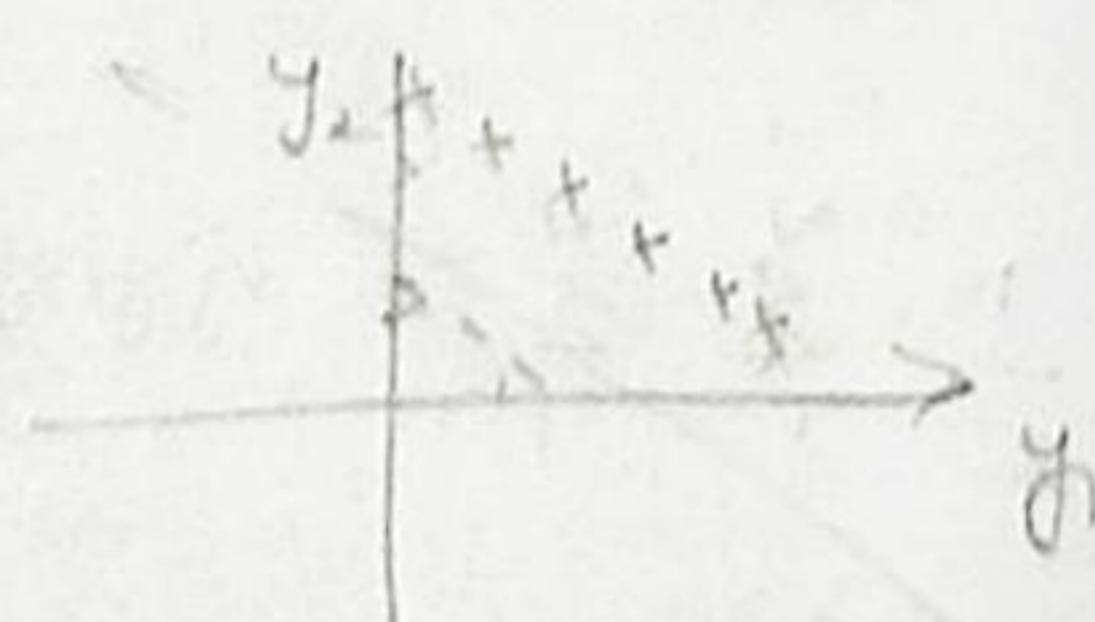
But a circle can do it beautifully.

Do a feature transform

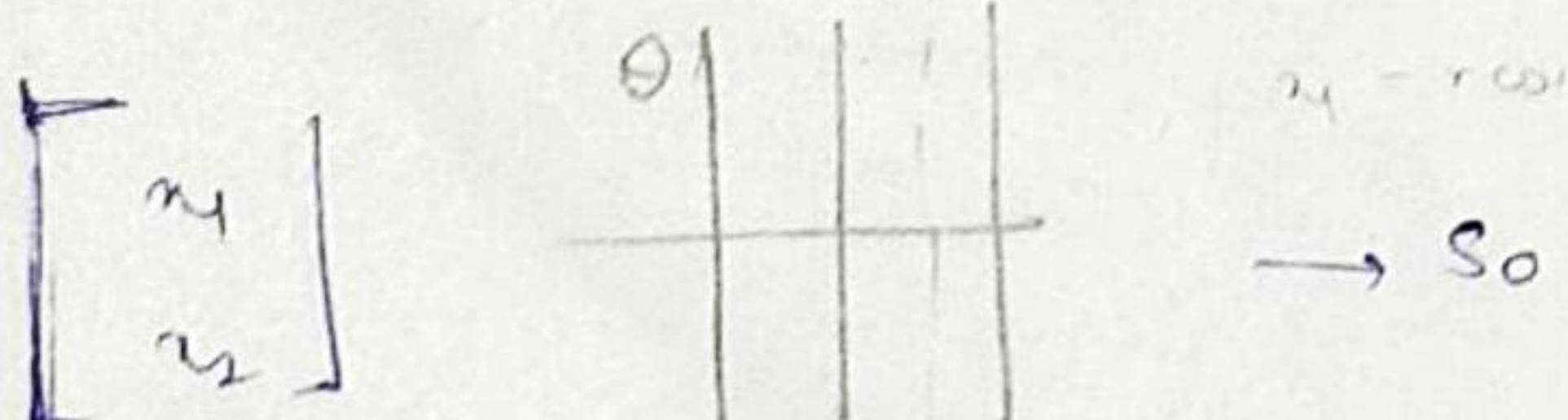
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} x_1^2 \\ x_2^2 \end{bmatrix}$$



Toy example



- If we knew how to transform, kernels can do stuff, ~~linear~~ classifiers can never do



→ So only r can separate

Separable

- we would like to map 2D space to  $\text{d}^N$  spaces

$$\bar{P} = \begin{bmatrix} P_1 \\ P_2 \end{bmatrix}, \quad \bar{q} = \begin{bmatrix} q_1 \\ q_2 \end{bmatrix}$$

$$P^T Q = P^T Q_1 + P^T Q_2$$

$$\text{Apply } \phi(P) \Rightarrow \begin{bmatrix} p_1^2 \\ p_2^2 \\ \sqrt{p_1 p_2} \end{bmatrix};$$

$$\phi(Q) = \begin{bmatrix} q_1^2 \\ q_2^2 \\ \sqrt{2}q_1q_2 \end{bmatrix}$$

$$\phi(P) \cdot \phi(O) = \begin{bmatrix} \sqrt{2}q_1q_2 \\ q_1^2 - q_2^2 \end{bmatrix} = P_1^2 q_1^2 + P_2^2 q_2^2 + 2P_1 P_2 q_1 q_2$$

$$K = (P^T Q)^2 \rightarrow \text{Homogeneous Polynomial Kernel}$$

$\phi(P) \cdot \phi(Q)$  is a polynomial.  
 we took  $P^T Q$  and squared it  
 $\Rightarrow$  we got dot product  $\phi(P) \cdot \phi(Q)$   
 without going to  $3D$ .

$$K(pq) = (p^T Q)^2$$

Kernel function

$$\kappa(pq) = \phi(p)^T \phi(q)$$

Find  $\phi$  so that  $k(pq) = (p^T q + 1)^2$  Inhomogeneous Poly.

$$(P_1 q_1 + P_2 q_2 + 1)^2 = \begin{bmatrix} P_1 \\ P_2 \end{bmatrix} \begin{bmatrix} q_1 \\ q_2 \end{bmatrix}$$

$$= \begin{bmatrix} p_1^2 \\ \alpha_1 p_2 \\ \beta_1 p_1 p_2 \\ \beta_2 p_1 \\ p_{p_2} \end{bmatrix} \quad \begin{bmatrix} q_1^2, \\ q_2, \\ s_2 q_2 q_1, \\ s_2 q_1, \\ 1 \end{bmatrix} \quad (\times 1) \circ B \circ$$

$e^{-\gamma x^T y}$ ,  $\tanh(x^T y) \rightarrow \underline{\infty D}$  [we generally don't care what these dimensions are]

$$\phi \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} = \begin{bmatrix} p_1^2 \\ p_2^2 \\ p_1 p_2 \\ p_2 - p_1 \end{bmatrix}$$

$$p_1 q_1 + p_2 q_2$$

$$\begin{aligned} & p_1^2 q_1^2 \\ & p_2^2 q_2^2 \\ & p_1^2 p_2^2 q_1^2 q_2^2 \\ & p_1^2 p_2^2 q_1^2 q_2^2 \end{aligned} = \frac{(p_1 q_1 + p_2 q_2)^2}{(p_1 q_1 + p_2 q_2)} = \frac{(p_1 q_1 + p_2 q_2)^2}{(p_1 q_1 + p_2 q_2)^2} = (p_1 q_1 + p_2 q_2)^2$$

so many  $\phi$  same kernel [Faisak!]

Kernel

if  $K$  is PSD,  $\exists \phi$

$$K = \left[ k(x_i, x_j) \right]_{N \times N}$$

- RBF: A Kernel that transforms to  $\infty$  space, can be used for most problems

Kernel need not be a dot product vector

$$K = K_1() + K_2() = \sum_{i=1}^n \alpha_i K_i()$$

Curse of dimensionality:-

It is bad, to use more dimensions

But Kernels use only dot prod ||, so no problem

$$\sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{K(x_i, x_j)}_{\phi(x_i)^T \phi(x_j)}$$

$$\text{sign} \left( \sum_{i=1}^n \alpha_i y_i K(x_i, x_j) + b \right)$$

$$\text{Sign}(\omega^T x)$$

$$\begin{aligned} & \omega_1 = 0 \\ & \omega_2 = 0 \\ & b = \omega_0 \end{aligned}$$

$$\text{sign}(\omega^T x + b)$$

• while testing  $(d_i, \vec{x}_i)$  are to be carried  
i.e SV

If we are going to calculate linear sum  
of  $d_i, \vec{x}_i$  data & store only  $\omega$ .  
we can remove

• If we know a feature map,