

## Lecture 7

\* Posing the Optimal margin classifier  
problem:

$$\begin{aligned} & \text{max } J \\ & \text{s.t. } \frac{y^{(i)}(w^T \cdot x^{(i)} + b)}{\|w\|} \geq \gamma \quad \text{"every example has at least } \gamma \text{ margin off axis"} \end{aligned}$$

maximize the  
want each  
ge one fine  
margin

~~$$\begin{aligned} y &= g(w^T x + b) \\ &= g\left(\frac{1}{\gamma} \frac{w^T x + b}{\|w\|}\right) \Rightarrow \end{aligned}$$~~

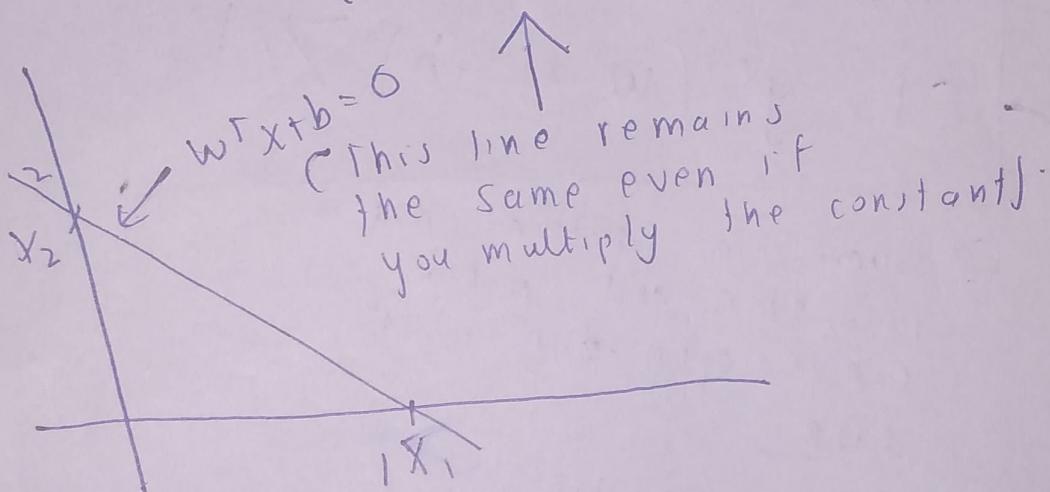
(65)

- ⑥ Example of a classifier

$$= g\left(\begin{bmatrix} 2 \\ 1 \end{bmatrix}^T x + b - 2\right) \quad \{ g(w^T x + b) \}$$

You can scale up the values of w and b, in this case by 10.

$$= g\left(\begin{bmatrix} 20 \\ 10 \end{bmatrix}^T x - 20\right)$$



- ⑦ You can multiply the equation  $g(w^T x + b)$  with any constant without changing the decision boundary.

⑧ Now returning back:-  
→ lets say you choose to normalize the equation by  $\|w\| = \frac{1}{J}$

then:-

$$\max \frac{1}{\|w\|}$$

s.t.  $y^{(i)}(w^T x^{(i)} + b) \geq 1$

(66)

$$\text{Max} \min \frac{1}{2} \|w\|^2$$

$$y^{(i)}(w^T x^{(i)} + b) \geq 1$$

④ Suppose  $w = \sum_{i=1}^m \alpha_i y_i^{(i)} x^{(i)}$  ( $w$  can

be written as a linear combination  
of training examples).

⑤ Why? (Reasonable Assumption above).

→ Intuition #1:

⑥ For logistic Regression

$$\theta_0 = 0$$

Gradient descent:

$$\theta_j = \theta_j - \alpha (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

→  $\theta$  is a linear combination  
of your training example

for logistic regression.

(Proof by induction)

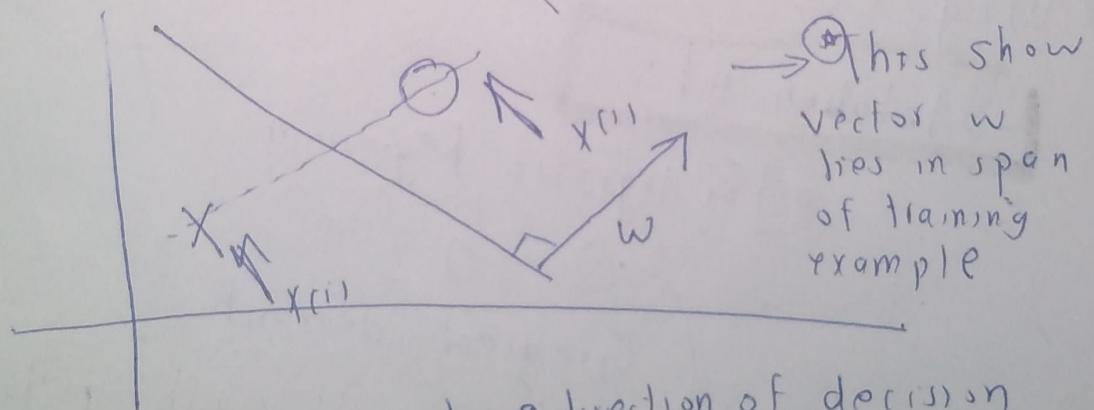
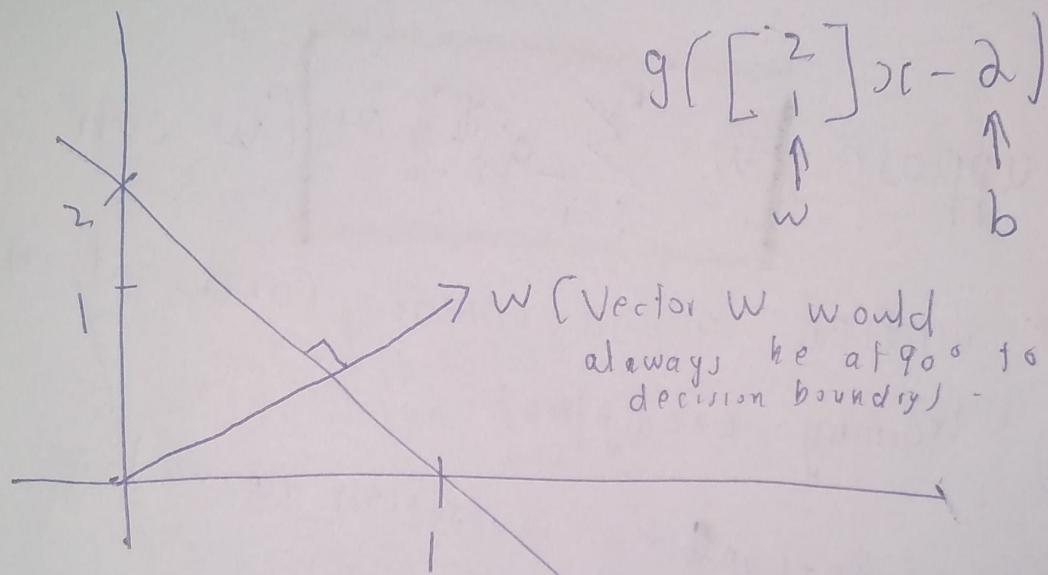
• (True for Batch Gradient  
Descent).

⑦ You can optimize SVM by gradient  
descent and see like logistic regression

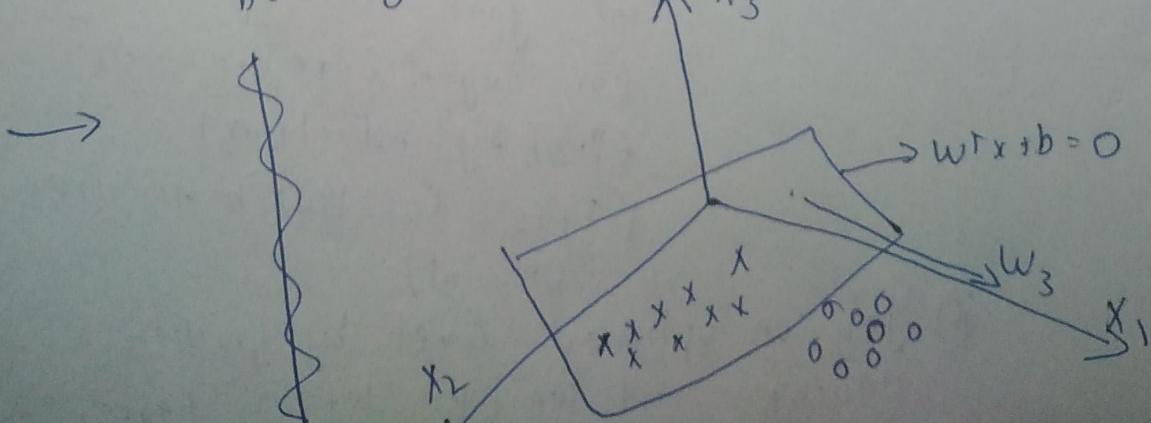
(67)

Parameter are linear combination of training examples.

→ Intuition 2



④ W set the positive direction of decision boundary while b shifts the decision boundary



(68)

→ The above diagrams shows that our training examples lies in three dimensional space but each training example third dimension is zero. The plane has been shown and vector normal to plane would be  $w$  which would also have third dimension as zero like training example which shows  $w$  upon training examples.

$$\rightarrow w = \sum_{i=1}^m \alpha^{(i)} y^{(i)} x^{(i)} \quad (\text{Plugging in below equations})$$

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y^{(i)} (w^T x^{(i)} + b) \geq 1 \quad i=1, \dots, m$$

$$\rightarrow \min_{w,b} \frac{1}{2} \left( \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T \left( \sum_{j=1}^m \alpha_j y^{(j)} x^{(j)} \right)$$

$$= \min \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \underbrace{x^{(i)T} x^{(j)}}_{\langle x^{(i)}, x^{(j)} \rangle}$$

$\langle x, z \rangle = x^T z$  is inner product

(69)

\*)

$$\text{s.t. } y^{(i)} \left( \sum_j \alpha_j y^{(j)} \langle x^{(j)}, x^{(i)} \rangle + b \right) \geq 1.$$

$$\text{s.t. } y^{(i)} \left( \sum_j \alpha_j y^{(j)} \langle x^{(j)}, x^{(i)} \rangle + b \right) \geq 1.$$

→ So Our updated equation are:

$$\min \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y^{(i)} y^{(j)} \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t. } y^{(i)} \left( \sum_j \alpha_j y^{(j)} \langle x^{(j)}, x^{(i)} \rangle + b \right) \geq 1.$$

- Feature vectors only appear in the inner product but you can compute inner product without actually mapping feature vector to higher dimension.

\*)

(10)

You further simplify the above equation  
to the following (Simplification is complicated)

$$\text{max} \sum_{i=1}^n \frac{1}{2} \sum_i \sum_j y_i y_j \alpha_i \alpha_j \langle x^i, x^j \rangle$$

$$\text{st } \alpha_i \geq 0 \quad \text{"Dual optimization problem"}$$

$$\sum_i y_i \alpha_i = 0$$

1) Solve for  $\alpha_i$ 's,  $b$

2) To make prediction compute

$$h_{w,b}(x) = g(w^T x + b)$$

$$= g\left(\underbrace{\left(\sum_i \alpha_i y_i x^i\right)^T}_{\tilde{w}} x + b\right)$$

$$= g\left(\sum_i \alpha_i y^i \langle x^i, x \rangle + b\right)$$

↳ You can now make prediction  
any way you want

(71)

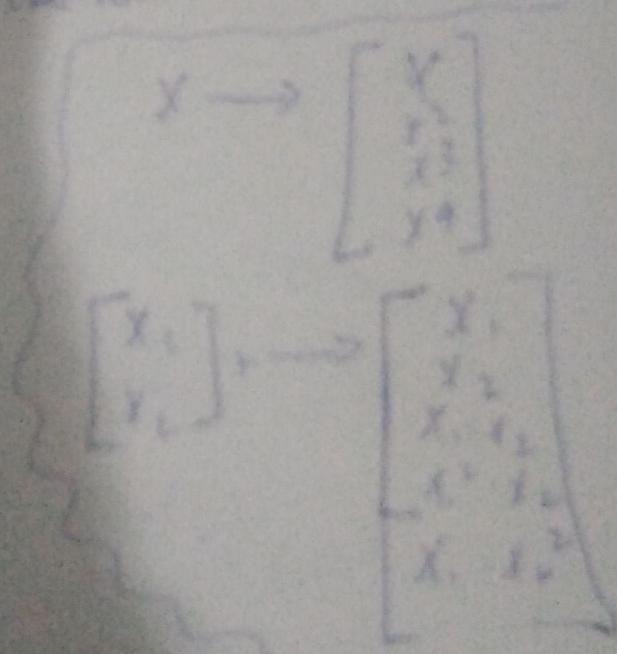
# Kernels Task

- 1) Write algorithm in terms of  $\langle \phi(x), \phi(z) \rangle$  or  $\langle x, z \rangle$
- 2) Let there be some mapping from

$x \mapsto \phi(x)$  - higher dimensional  
map  
higher  
dimension

2) Find way to compute  
 $\langle \phi(x), \phi(z) \rangle = \phi(x)^T \phi(z)$

↓  
 Kernel function (allow to  
 compute dot product even  
 if  $\phi(x), \phi(z)$  are in  
 very high dimensions)



(72)

4) Replace  $\langle x, z \rangle$  in algorithm  
with  $K(x, z)$ .

① Examples of kernel

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

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

$$\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix}$$

$$\phi(x) \in \mathbb{R}^{n^2}$$

$$\phi(z) = \begin{bmatrix} z_1 z_1 \\ z_2 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_3 \\ z_1 z_1 \\ z_2 z_1 \\ z_2 z_3 \end{bmatrix}$$

$$\phi(z) \in \mathbb{R}^{n^2}$$

→ Now because there are  $n^2$  elements you need  $O(n^2)$  time to compute  $\phi(x)$ , or  $\phi(x)^T \phi(z)$  explicitly.

$$\rightarrow K(x, z) = \phi(x)^T \phi(z)$$

$$= (x^T z)^2 \quad \begin{array}{l} \text{This is an} \\ O(n) \text{ time} \\ \text{operation.} \end{array}$$

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

Let me proof this statement

$$= \phi^T(x) \phi(z) = (x^T z)^2 = \left( \sum_{i=1}^n x_i z_i \right) \left( \sum_{j=1}^n x_j z_j \right)$$

$$= \sum_{i=1}^n \sum_{j=1}^n x_i z_i x_j z_j$$

$$= \sum_{i=1}^n \sum_{j=1}^n (x_i x_j) (z_i z_j) \rightarrow$$

$$\phi(x)$$

Poly nomial  
pairs

$$\phi(z)$$

Poly nomial  
pairs

Dot product  
of the vector  
which we mapped  
into higher dimension  
by polynomial pairs

(79)

→ We have been to compute the high dimension dot product  $\phi(x)^T \phi(z)$  by remaining in lower dimension through the kernel trick or shown by the example (Computationally Efficient).

## 5 → Other Examples of Kernel

2. ①  $K(x, z) = (x^T z + c)^2$

Cf IR  
This correspond to modifying your features as follow

$$\phi(x) = \begin{bmatrix} x_1 y_1 \\ x_1 y_2 \\ x_1 x_3 \\ x_2 y_1 \\ x_2 y_2 \\ x_2 x_3 \\ x_3 y_1 \\ x_3 y_2 \\ x_3 x_3 \\ \bar{x}_1 \bar{x}_1 \\ \bar{x}_1 \bar{x}_2 \\ \bar{x}_1 \bar{x}_3 \\ \bar{x}_2 \bar{x}_1 \\ \bar{x}_2 \bar{x}_2 \\ \bar{x}_2 \bar{x}_3 \\ \bar{x}_3 \bar{x}_1 \\ \bar{x}_3 \bar{x}_2 \\ \bar{x}_3 \bar{x}_3 \end{bmatrix}$$

→ Same change  
to  $\phi(z)$ .

75

→ So the simple  $K(x, z)$  can compute the dot product of transformed  $\phi(x) \cdot \phi(z)$  without any mapping

$$\rightarrow K(x, z) = (\phi^T z + c)$$

$$\rightarrow K(x, z) = (x^T z + c)^d \rightarrow O(n) \text{ time computation}$$

This correspond to taking dot product of

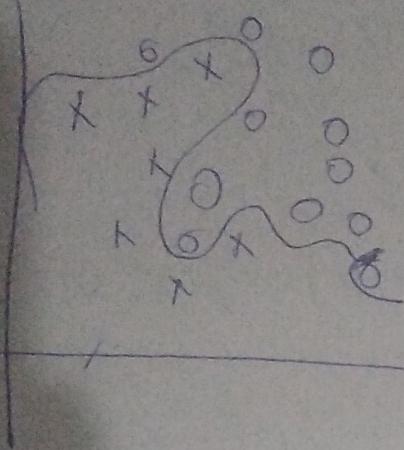
④  $\phi(x) + \text{has all } \binom{n+d}{d} \text{ features}$

of monomial upto order  $d$ .

→ So **SVM** is the optimal margin classifier and applying kernel trick to

11

$\rightarrow$  SVM = Optimal Margin + Trick



→ This is a linear boundary learned by SVM in high-dimensional space but when you project back to original space it becomes non-linear.

## ④ How to make kernels.

→ If  $x, z$  are "similar"  $K(x, z) = \phi(x)\phi(z)$

is large

• If  $x, z$  are "dissimilar",  $K(x, z)$  is small

$$\rightarrow K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

→ You can see that kernel are basically similarity mapping function so  $K(x, z)$  basically satisfies this idea where if  $x$  and  $z$  are similar its value will be high or else vice-versa. (Satisfies the criterion).

(21) You can use the Kernel function such as the one suggested above only if

there exist  $\phi$  s.t  $K(x, z) = \phi(x)^T \phi(z)$

Thus  $K(x, x) = \phi(x)^T \phi(x) \geq 0$

↳ Dot product  
is also always  
greater than zero

→ Proof to outline when

Kernel is valid

det  $\{x^{(1)}, \dots, x^{(d)}\}$  any d points

det  $K_{ER}^{d \times d} \rightarrow$  "Kernel matrix"

$K_{ij} = K(x^{(i)}, x^{(j)}) \rightarrow$  Kernel applied  
to two of those  
points. (Every pair  
of point applied with  
Kernel and put into  
the matrix)

78

④ Given any vector  $\mathbf{z}$ .

$$\mathbf{z}^T \mathbf{K} \mathbf{z} = \sum_i \sum_j z_i K_{ij} z_j$$

$$= \sum_i \sum_j z_i (\phi(x^{(n)})^T) (\phi(x^{(n)}) z_j)$$

~~$\leftarrow \text{if } K \neq 0$~~

$$= \sum_i \sum_j z_i \sum_k (\phi(x^{(n)})_k (\phi(x^{(n)})_k)) z_j$$

$$= \sum_k \left( \sum_i z_i (\phi(x^{(n)})_k) \right) \left( \phi(x^{(n)})_k z_j \right)$$

$$= \sum_k \left( \sum_i z_i (\phi(x^{(n)})_k) \right)^2 \geq 0$$

which proves  $K \geq 0$  ( $K$  is positive

semi-definite matrix).

79

$\Rightarrow$  Theorem Mercer:  $K$  is a valid

Kernel function ( $\in \mathbb{F}$ ) if & st

$K(x, z) = \phi(x)^T \phi(z)$  if and only  
if for any  $d$  points  $\{x^{(1)}, x^{(d)}\}$   
the corresponding Kernel matrix  
 $K \geq 0$  (Positive semi-definite).

$\hookrightarrow$  The theorem also stands  
in a reverse way i.e. if  $K \geq 0$  then

$K(x, z) = \phi^T(x) \phi(z)$  is a valid

Kernel.

$$\rightarrow K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

$\downarrow$   
Gaussian  
Kernel

\* Different types of Kernel - (80)

### Linear Kernel

$$K(x, z) = x^T z$$

$$\phi(x) = x$$

### Gaussian Kernel

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right) \rightarrow \text{This also uses all polynomial features.}$$

$$\phi(x) \in \mathbb{R}^\infty \text{ (Infinite dimension)}$$

### Polynomial Kernel

$$K(x, z) = (x^T z + c)^d \quad \text{($n \times d$) features}$$

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

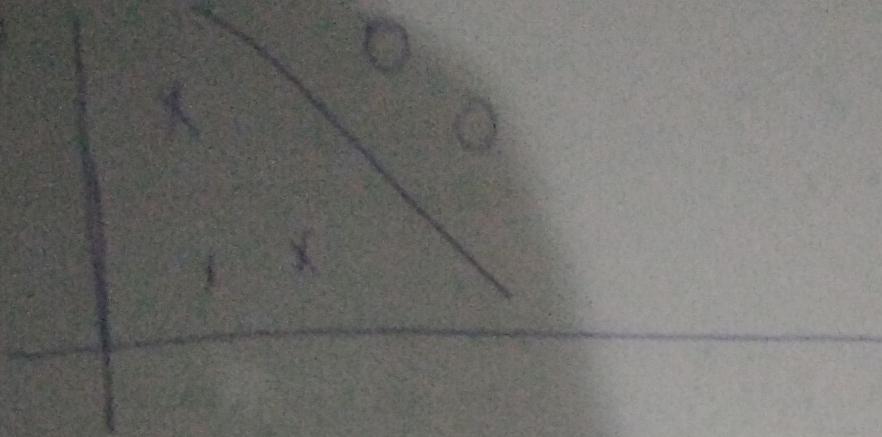
④ Kernel Trick is more general than

④ Kernel Trick

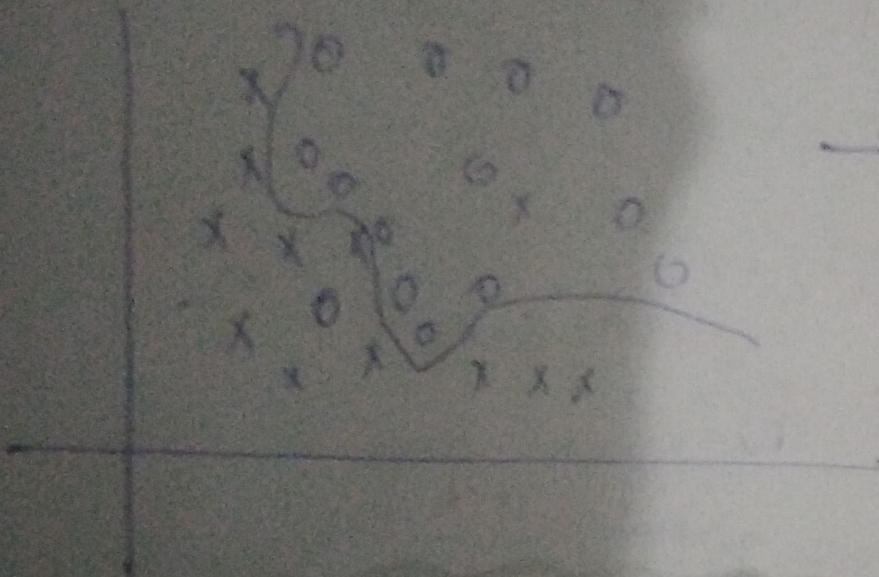
SVM:-  
→ You can apply Kernel trick algorithm if it can be rewritten in terms of  $\langle \phi(x^{(1)}), \phi(x^{(2)}) \rangle$   
and then replace it with Kernel function.

All the discriminative algorithms learned so far can take advantage of Kernel trick.

(8)



③ If Your dataset is noisy:



→ You don't want a very strict complicated decision boundary but rather a more soften one so it generalizes better (AKA You don't want a zero training error).

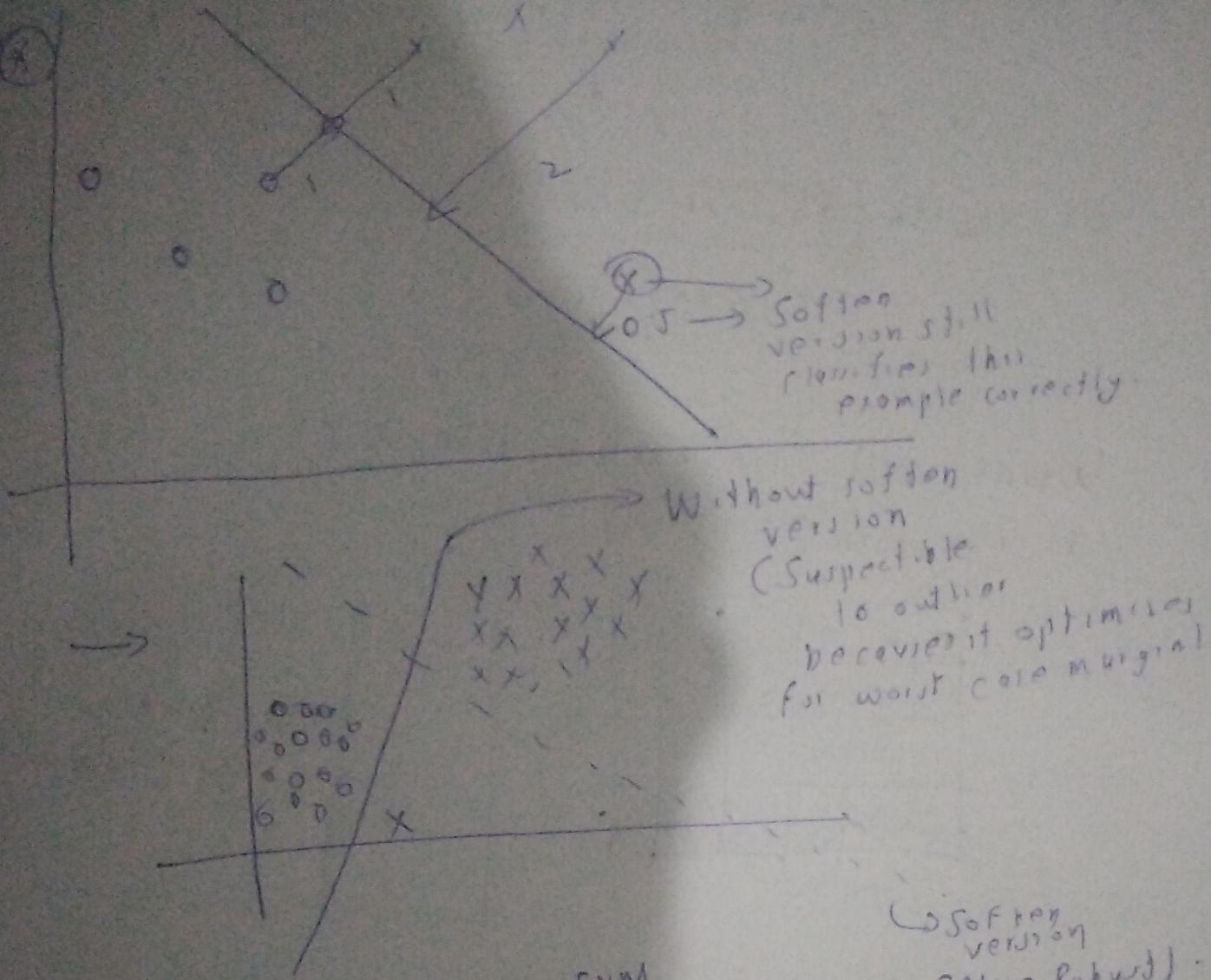
Q4 norm soft margin SVM

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i = 1, \dots, m.$$

soft margin

7/10



→ L, norm soft margin SVM  
Simplifies to the following:

$$\text{min} \sum_{i=1}^m d_i - \sum_{i=1}^m \sum_{j=1}^m y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s.t. } \sum_{i=1}^m y_i \alpha_i = 0$$

$$0 \leq \alpha_i \leq C$$

thus the change.

From previous optimization  
of hard SVM.

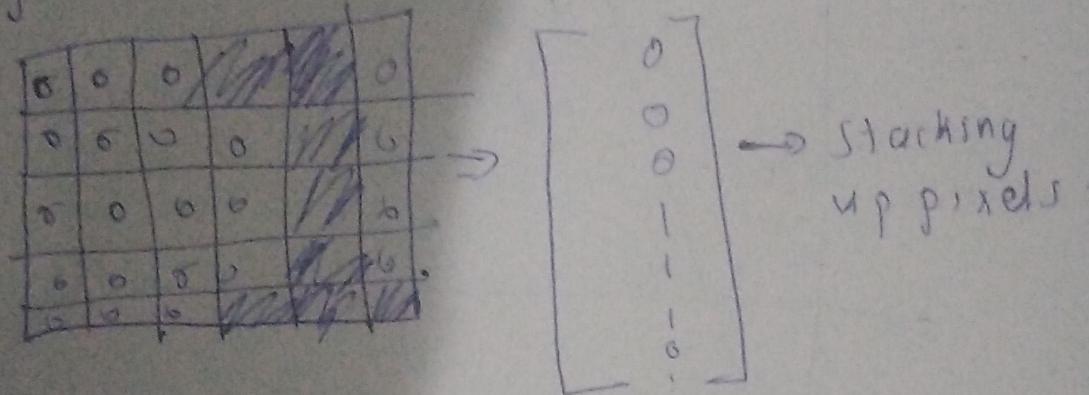
## ⑧ SVM with Polynomial Kernel

$$K(x, z) = (x \cdot z)^d$$

or with Gaussian Kernel

$$= \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right).$$

Works well with hand written digit recognition (Classifier).



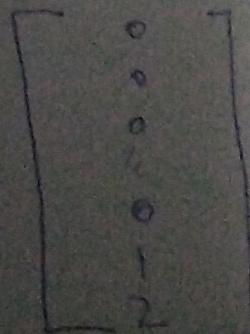
## ⑨ Protein Sequence classifier:

A...Z → Denote amino acid who combination define protein.

B A T T P A I B A J T T A U → Protein

$$\phi(x) = ?$$

but  
all four  
characters  
are different  
BAJT



Each feature represent how many times four continuous character appears in the protein.

$$\phi(x) = ?$$

→ 20<sup>9</sup> Dimensional feature.

(84)

→ write out instead of programming

$\phi(x)$  which is do a dimensional vector  
and consequently computationally

expensive to do so  $\Rightarrow$  You can compute  
kernel through dynamic programming re

$\phi(x)^T \phi(z) = K(x, z)$ . and then you use

support vector machine

↳ Example of innovative  
kernel.