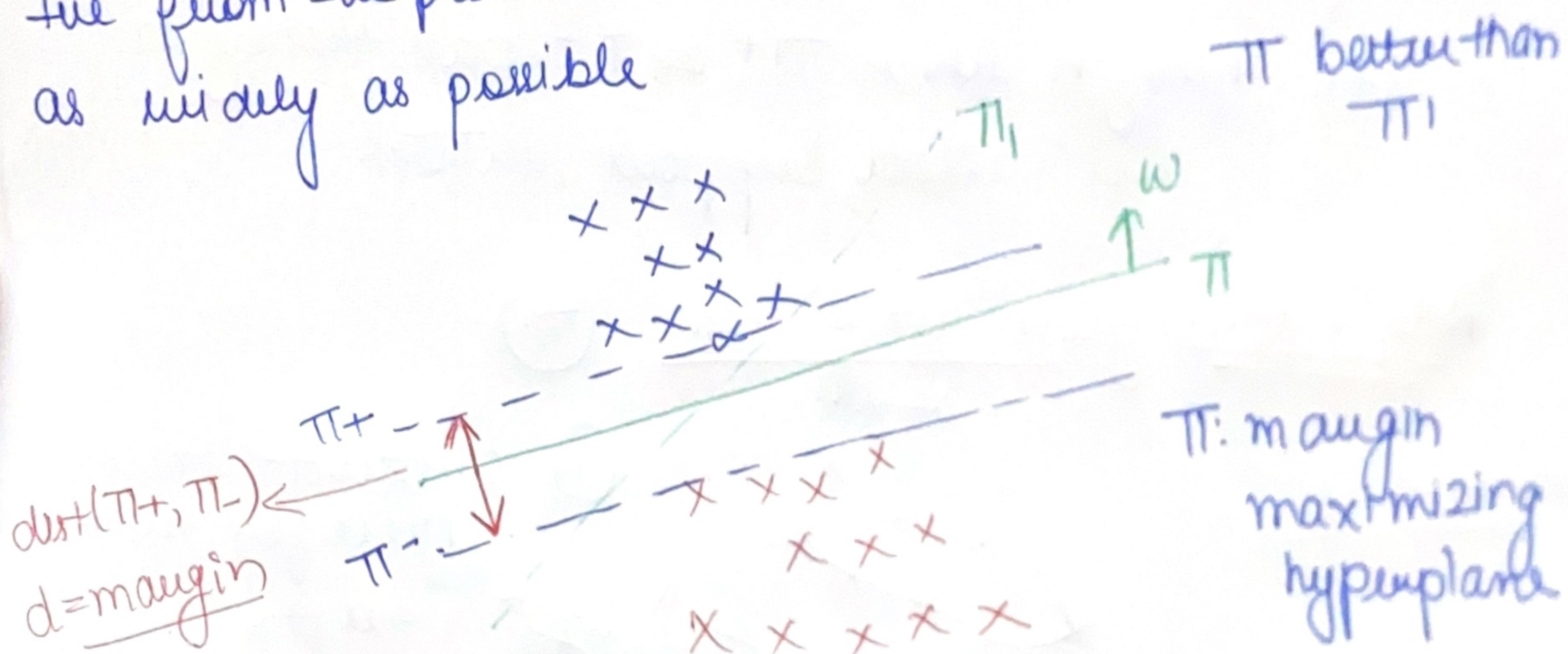
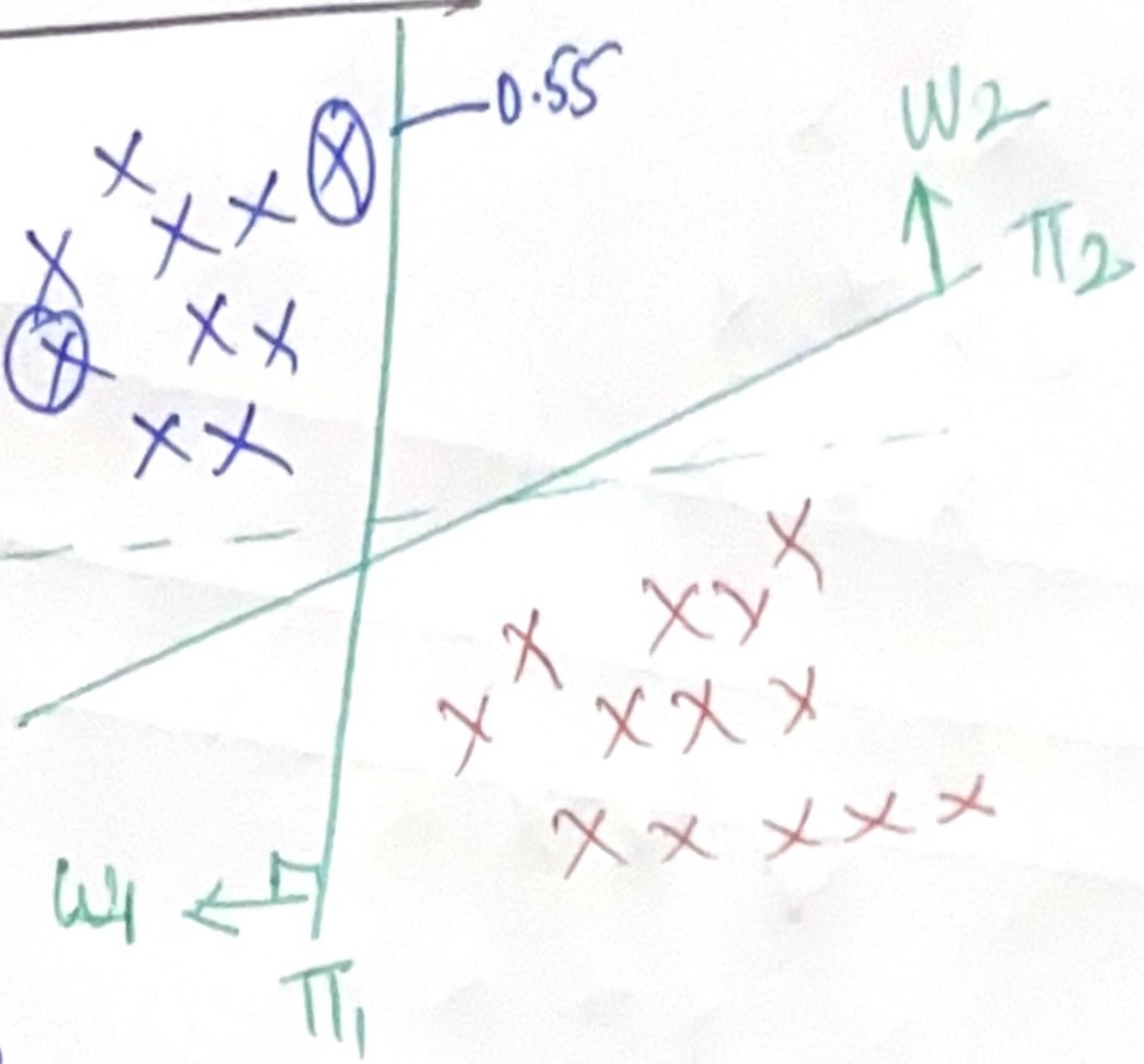


Support Vector Machine → classfn  
→ regression

### Geometric Intuition

many  $\pi$ s that separate the +ve from -ve

→ key idea of SVM:  
 $\pi$  that separates the +ve from -ve pts as widely as possible



$\pi_+$  is  $\parallel$  to  $\pi$

$\pi_-$  is  $\parallel$  to  $\pi$

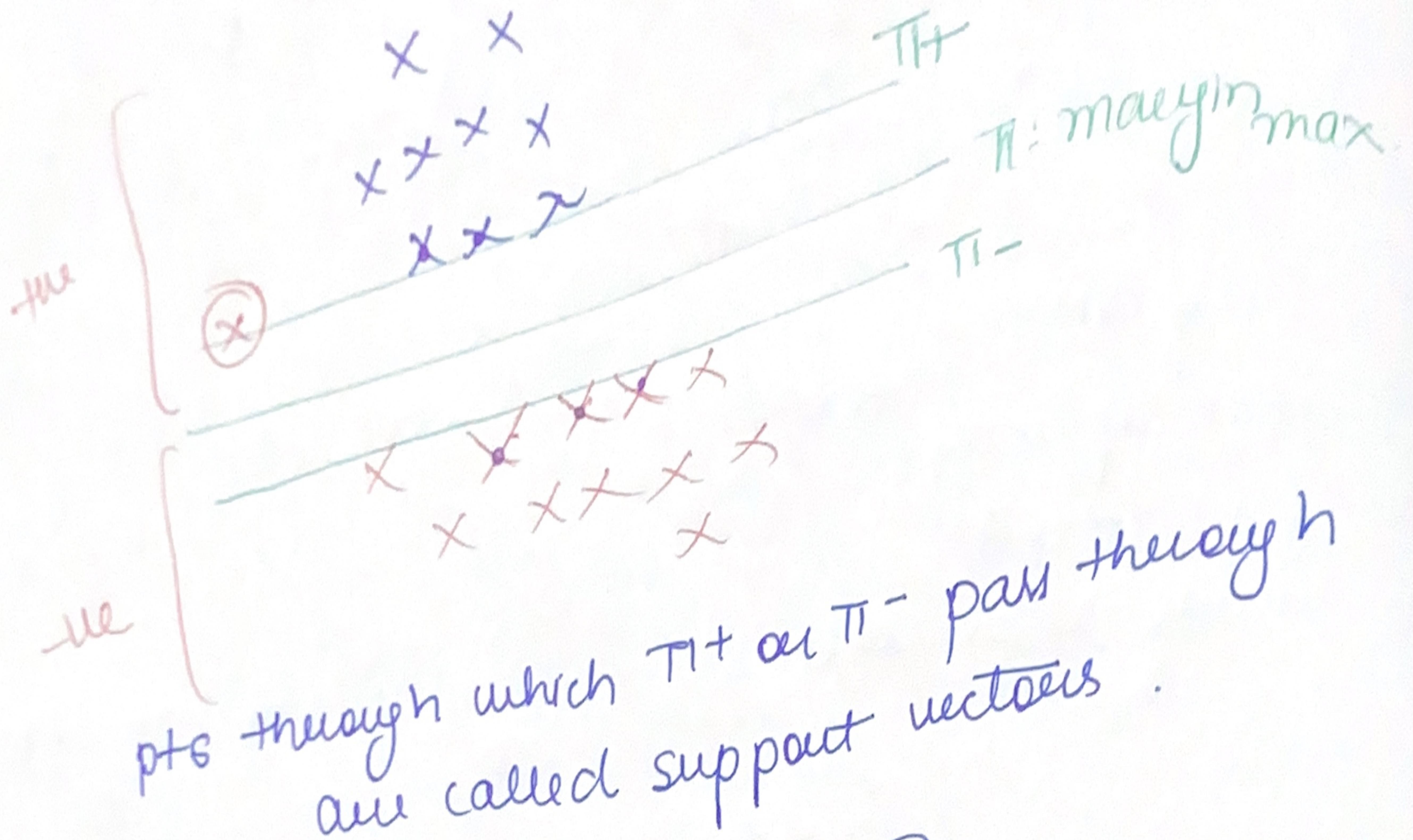
$\pi_+$  &  $\pi_-$  are  $\perp$  to  $\pi$

$\pi_+$  &  $\pi_-$  are  $\perp$  to each other

SVM: Try to find hyperplane which maximises margin  
 $\Rightarrow d_{\pi}(\pi_+, \pi_-)$

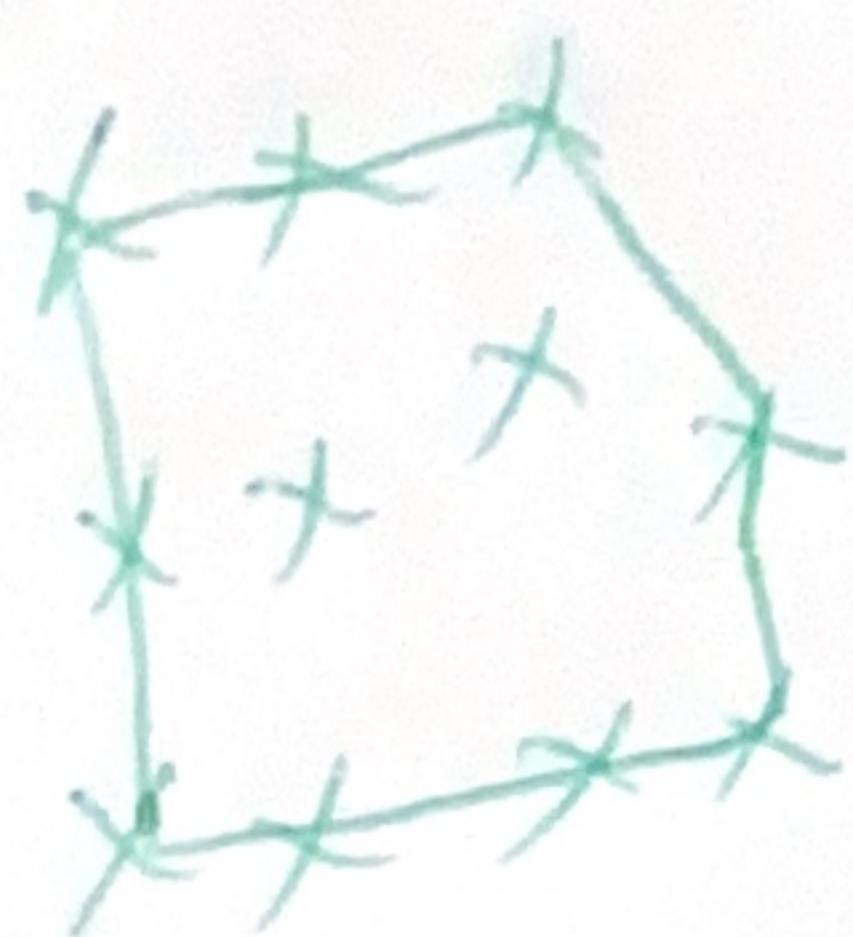
margin  $\uparrow$  generalization acc  $\uparrow$

# Support Vector

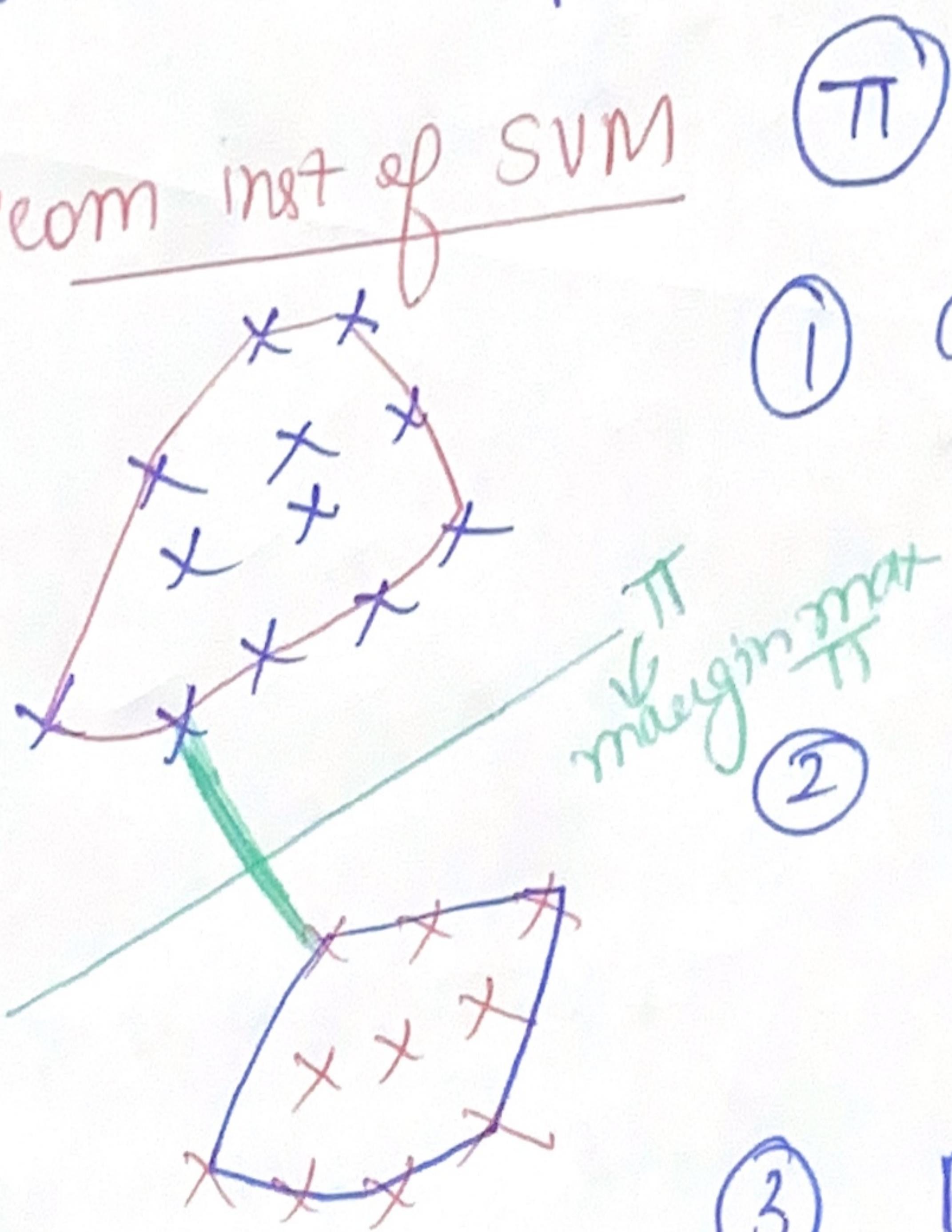


## Altundhile geom inst of SVM

### Convex-hull



Convex polygon



(1)

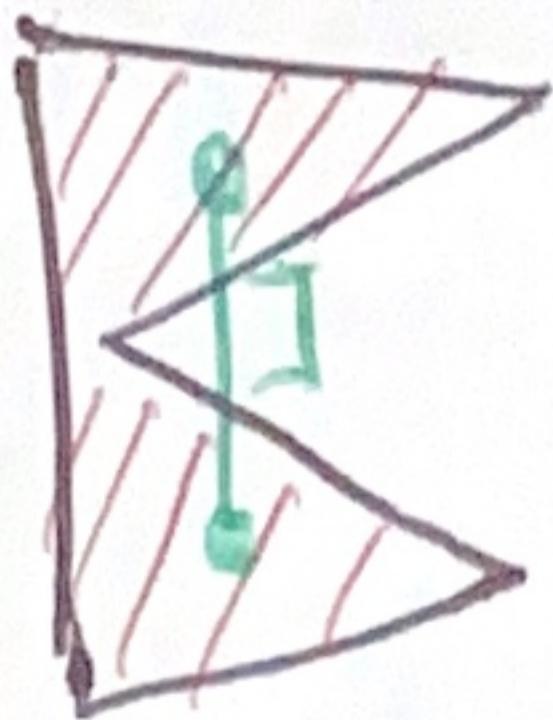
construct convex hull for  $n$  pts.

(2)

find shortest line connecting hulls

(3)

Bisect line



non  
convex  
polygon

## Mathematical formulation of SVM

$\Pi$ : margin maximization

$$\Pi: \mathbf{w}^T \mathbf{x} + b = 0$$

$$\text{if } \Pi^+: \mathbf{w}^T \mathbf{x}_i + b = 1$$

$$\Pi^-: \mathbf{w}^T \mathbf{x}_i + b = -1$$

$$\text{Margin: } d = \frac{2}{\|\mathbf{w}\|}$$

such that

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \text{ for all } \mathbf{x}_i$$

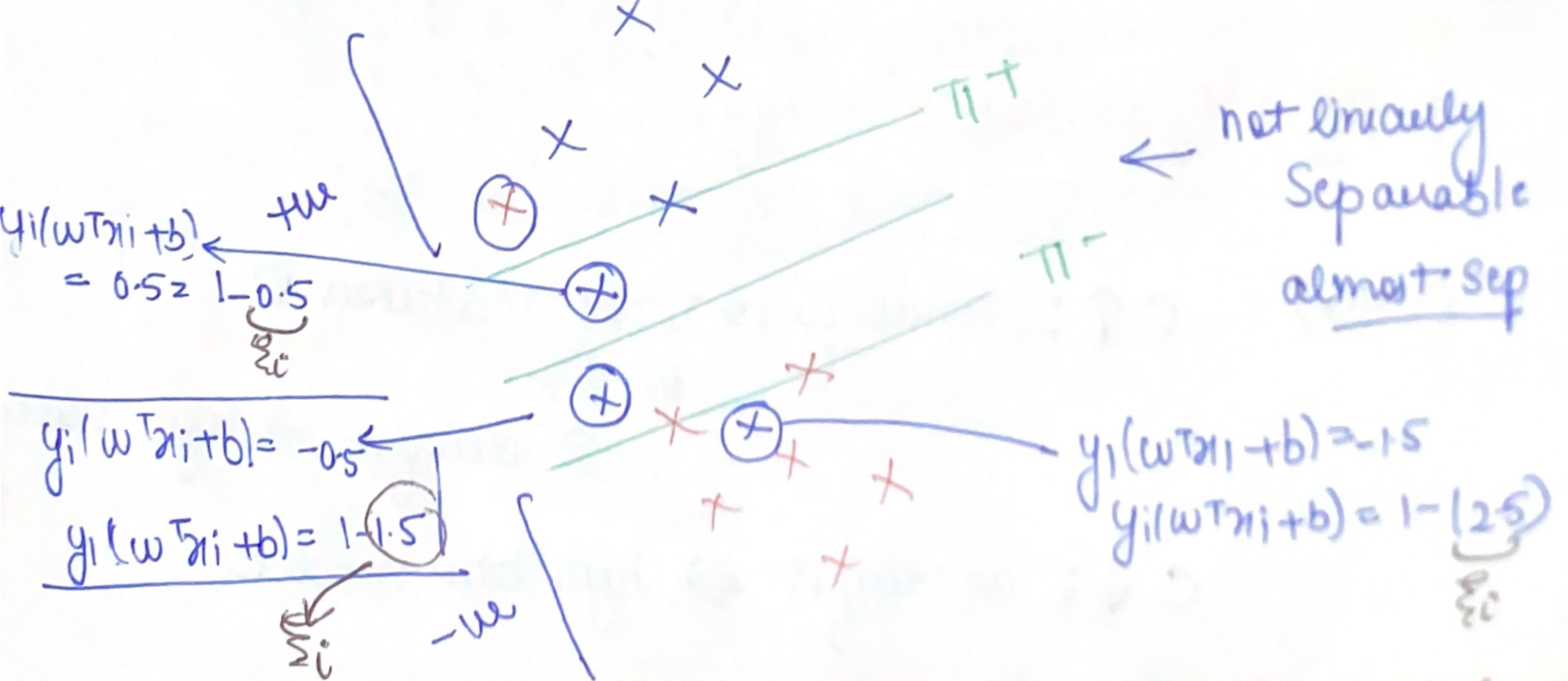
$$\mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\text{augmax}} \frac{2}{\|\mathbf{w}\|} = \text{margin.}$$

s.t.

$$\left[ \begin{array}{l} \mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\text{augmax}} \frac{2}{\|\mathbf{w}\|} \\ \text{s.t. } y_i, y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \end{array} \right]$$

constraint optimization problem of SVM

data is linearly separable and margin  $\Sigma_i$



$\xi_i \uparrow$ : pt  $x_i$  further away from correct  $\pi$

$$x_i \rightarrow \xi_i$$

$$\xi_i = 0 \text{ if } y_i(w^T x_i + b) \geq 1$$

(Q)

$\xi_i > 0$   $\pi$  is equal to some unit of dist away from correct hyperplane in incorrect direction

$$w^*, b^* = \underset{w, b}{\text{augmax}} \frac{2}{\|w\|} = \underset{w, b}{\text{augmin}} \frac{\|w\|}{2}$$

hyperpara -  $\max f(x) = \min \frac{1}{f'(x)}$

$$w^*, b^* = \underset{w, b}{\text{augmin}} \frac{\|w\|}{2} + C \cdot \frac{1}{n} \sum_{i=1}^n \xi_i \rightarrow \text{avg dist of misclassified pt from } \pi$$

$y_{\text{margin}}$

$$\begin{aligned} & \text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \\ & \xi_i \geq 0 \end{aligned} \quad \begin{array}{l} \text{correctly classif pt} \\ \xi_i = 0 \end{array}$$

Soft margin SVM

minimize errors = min misclassification



$$\min \sum \xi_i$$

$$\min_w (\text{logistic loss}) + \alpha (\text{reg})$$

$C(+ve)$   $C \uparrow$ ; tendency to make mistakes  $\downarrow$   
on  $D_T$

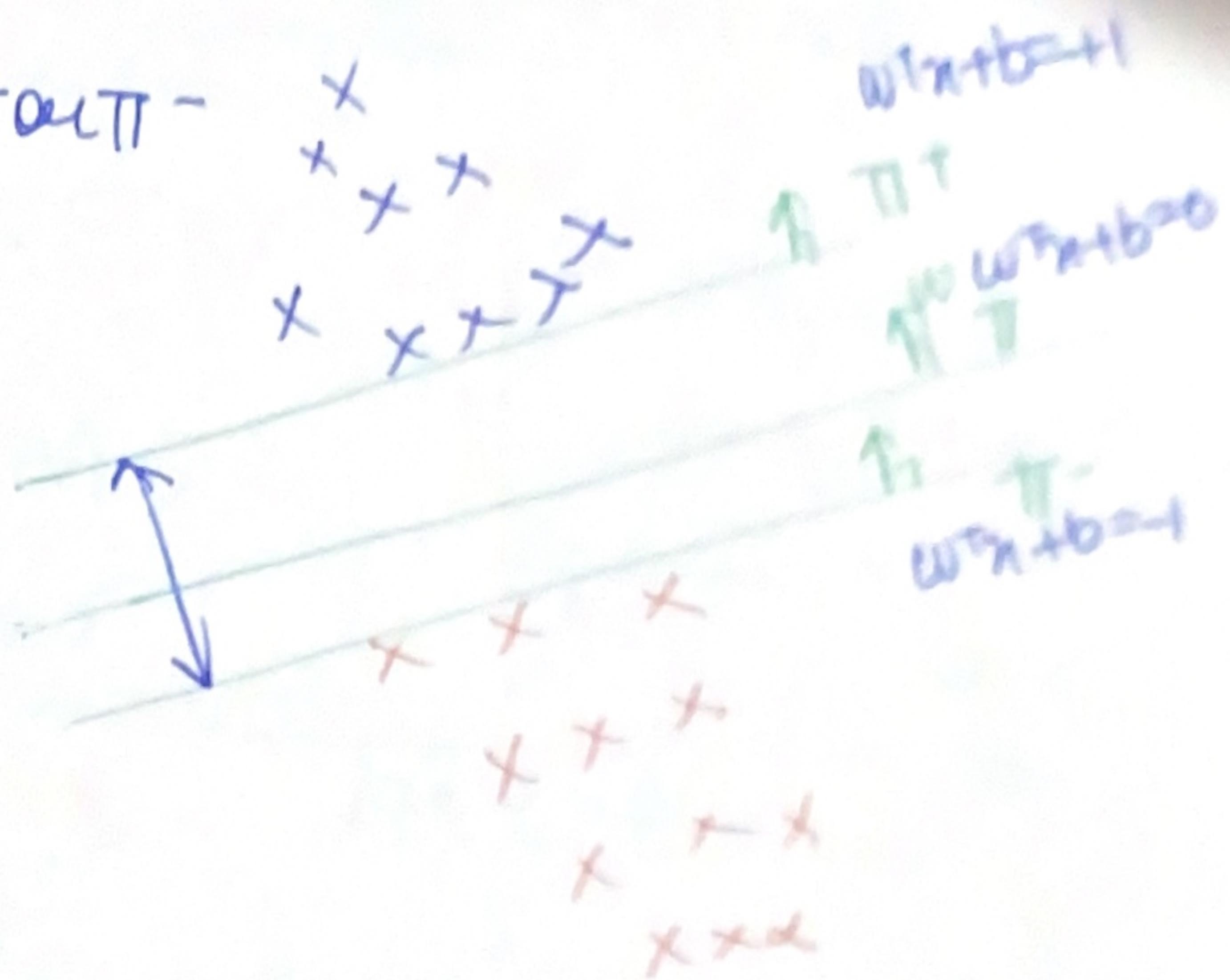
$\Rightarrow$  overfit  $\rightarrow$  high var

$C \downarrow$ ; underfit  $\Rightarrow$  high bias model

## SVM:

(Q) Why +8-1 on RHS of  $\pi^+ \text{or} \pi^-$

margin:  $\frac{2}{\|\omega\|}$



$$w^*, b^* = \underset{b}{\operatorname{arg\,max}} \frac{2}{\|\omega\|}$$

$\|\omega\| \neq 1$  (any vector)  
need not be unit vector

①  $\pi^+$ :  $w^T x + b = 1$

$$\pi^-: w^T x + b = -1 \quad K > 0$$

margin:  $\frac{2K}{\|\omega\|}$

$$K=4$$

$$\underset{w,b}{\operatorname{arg\,max}} \frac{2}{\|\omega\|} = \underset{w,b}{\operatorname{arg\,max}} \frac{2K}{\|\omega\|} = \frac{8}{\|\omega\|}$$

②  $w^T x + b = 1$

$$\left(\frac{w}{K}\right)^T x + \left(\frac{b}{K}\right) = 1 \quad w \perp \pi$$

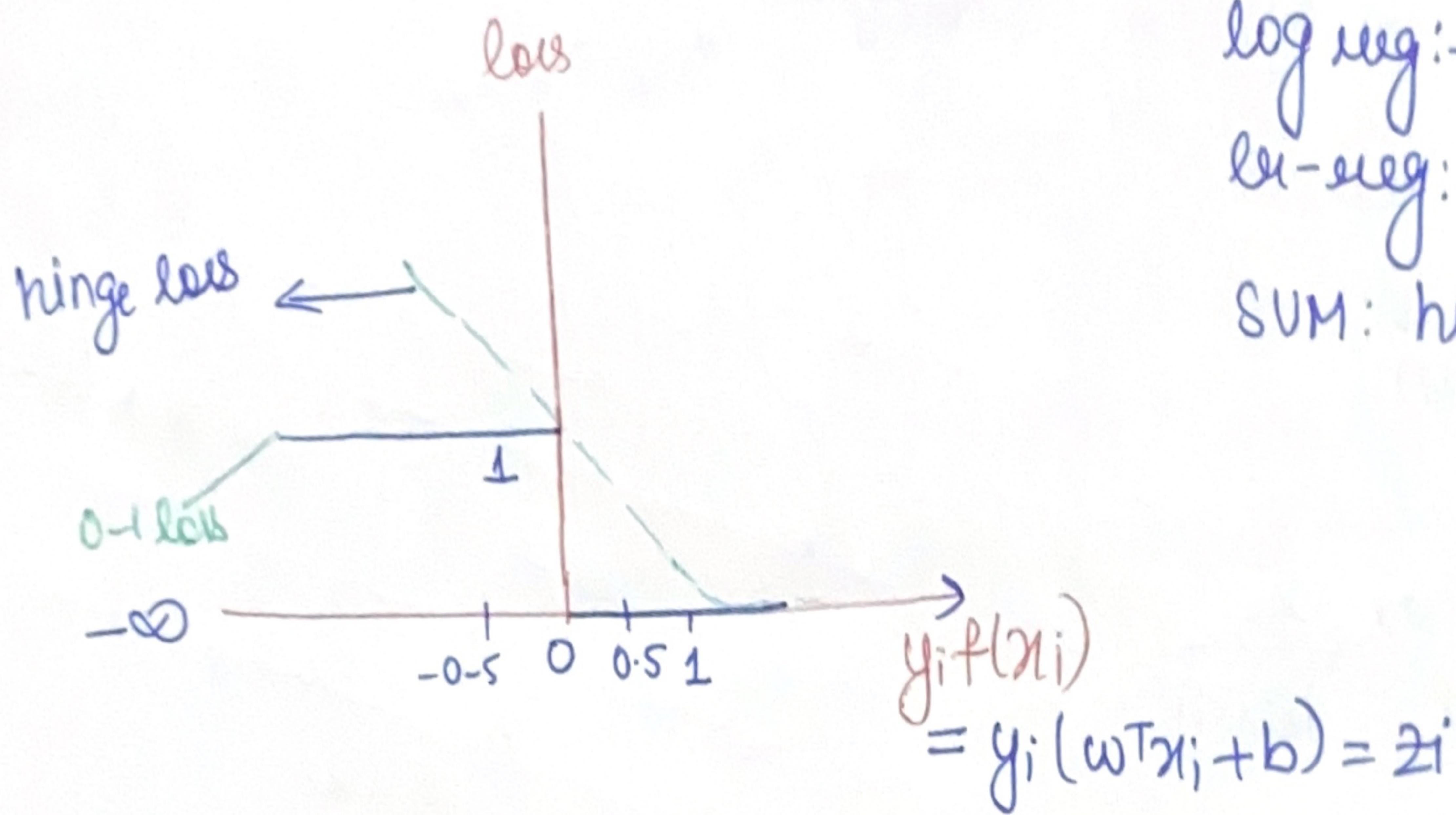
$$\|\omega\| \neq 1$$

$$(w^1)^T x + b^1 = 1$$

Raise +1 or -1

↑ Simplify math

## Loss minimization : Hinge-loss



$\begin{cases} z_i \geq 0: x_i \text{ is correctly classified} \\ z_i < 0: x_i \text{ is incorrectly classified} \end{cases}$

Hinge loss:  $z_i \geq 1$ ; hinge loss = 0

$z_i < 1$ ; hinge loss =  $1 - z_i$

$$\max(0, 1 - z_i)$$

Case 1:  $z_i \geq 1$ ;  $1 - z_i$  is -ve value  $\Rightarrow \max(0, 1 - z_i) = 0$

Case 2:  $z_i < 1$ ;  $1 - z_i > 0 \Rightarrow \max(0, 1 - z_i) = 1 - z_i$

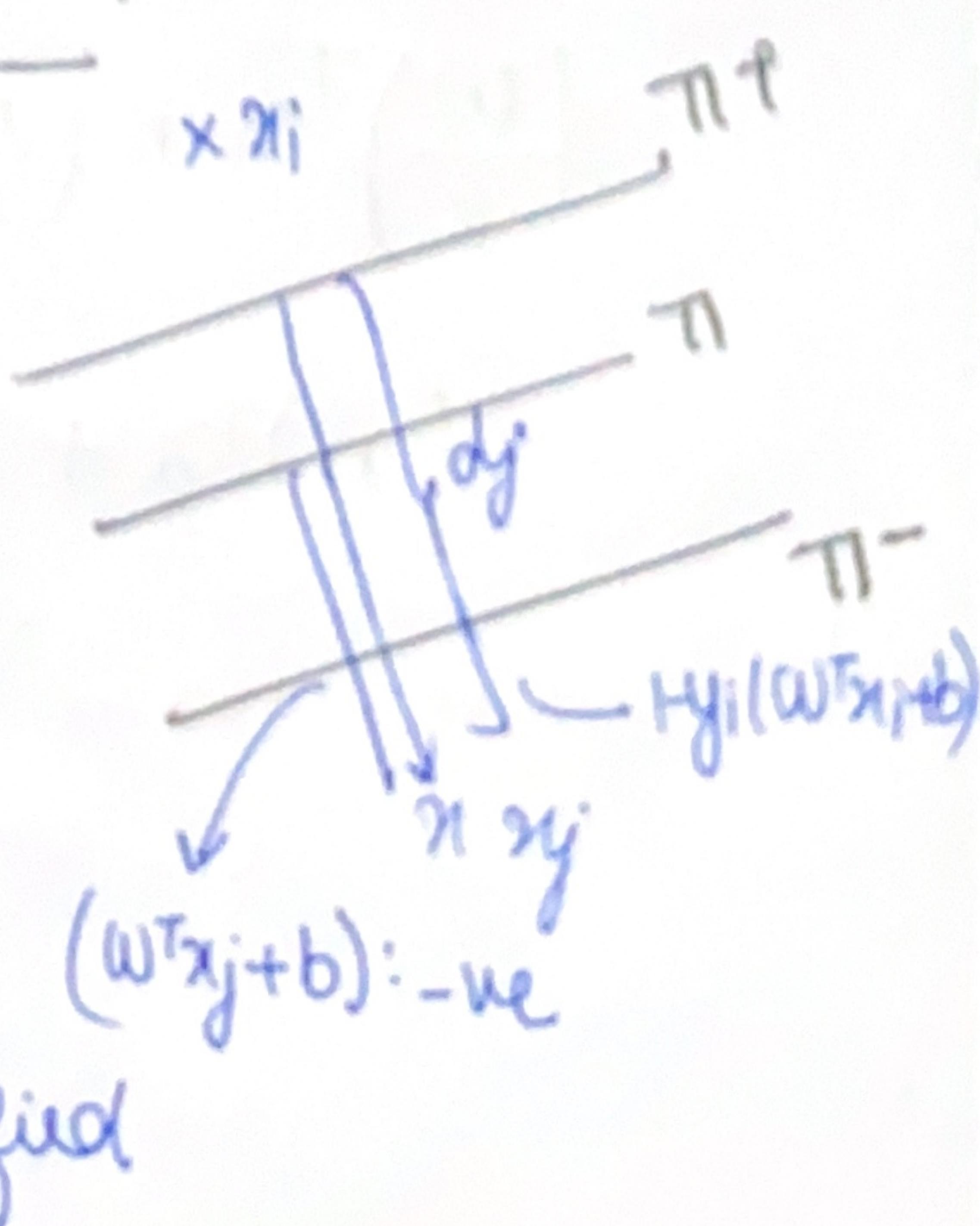
## Geometric formulation & learning

$$\xi_j = 0 \leftarrow x_j$$

$$d_j = 1 - y_j(\omega^T x_j + b) = 1 - z_j$$

$$\xi_j = \text{dist of } x_j \text{ to the } \pi^+ = d_j = \frac{1 - z_j}{\|\omega\|}$$

$\xi_j = 1 - z_j$  when  $x_j$  is misclassified



$$\max(0, 1 - \xi_i) = \xi_i$$

Soft SVM

$$\begin{aligned} & \min_{w,b} \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i \\ \text{s.t. } & 1 - y_i(w^T x_i + b) \geq \xi_i \\ & \xi_i \geq 0 \end{aligned}$$

$C \uparrow \Rightarrow$  overfit

$C \downarrow \Rightarrow$  underfit

~~hard SVM~~  
loss min

$$\begin{aligned} & \min_{w,b} \sum_{i=1}^n \max(0, 1 - y_i(w^T x_i + b)) + \\ & \quad \lambda \|w\|^2 \end{aligned}$$

$\lambda \uparrow \Rightarrow$  underfit       $\lambda \downarrow \Rightarrow$  overfit

$$\|w\| > 0 \Rightarrow \min \frac{\|w\|}{2} \text{ is same as } \|w\|^2$$

### Dual form of SVM

soft margin SVM

$$\begin{aligned} & \min_{w,b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t. } & y_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \\ & \xi_i \geq 0 \end{aligned}$$

$\xrightarrow{\text{Equivalent}}$  primal of SVM

$$\begin{aligned} & \max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t. } & \alpha_i \geq 0 \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

- ①  $x_i \rightarrow \alpha_i$
- ②  $\alpha_i$ 's only occur in form of  $x_i^T x_j$
- ③  $f(x_q) =$

Dual form

$$③ f(\mathbf{x}_q) = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^\top \mathbf{x}_q + b$$

SVs:  $\alpha_i > 0$

non SVs:  $\alpha_i = 0$

$f(\mathbf{x}_q)$ : only pts that matter are SVs

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j$$

Sim( $\mathbf{x}_i, \mathbf{x}_j$ )

$\text{s.t. } \alpha_i \geq 0 \rightarrow \alpha_i = 0 \text{ for SVs}$

$\sum_{i=1}^n \alpha_i y_i = 0 \text{ for nonSVs}$

$$\mathbf{x}_i^\top \mathbf{x}_j = \mathbf{x}_i \cdot \mathbf{x}_j = \text{Cosine sim}(\mathbf{x}_i, \mathbf{x}_j)$$

if  $\|\mathbf{x}_i\| = 1 ; \|\mathbf{x}_j\| = 1$

$$f(\mathbf{x}_q) = \sum_{i=1}^n \alpha_i y_i \underline{\mathbf{x}_i^\top \mathbf{x}_q} + b$$

$\rightarrow k(\mathbf{x}_i, \mathbf{x}_q)$

## Kernel Trick

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \rightarrow \text{Sim}(\mathbf{x}_i, \mathbf{x}_j)$$

s.t.  $\sum_{i=1}^n \alpha_i y_i = 0 ; \alpha_i \geq 0$   $\downarrow k(\mathbf{x}_i, \mathbf{x}_j)$

Kernel func

$$f(\mathbf{x}_q) = \sum_{i=1}^n \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}_q) + b$$

- The most impo idea in SVM is Kernel Trick
- Soft-SVM-hyperplanes ↪ log reg  
↓  
margin-max

for SVM;  $\mathbf{x}_i^T \mathbf{x}_j$        $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$   
Kernel SVM;  $k(\mathbf{x}_i, \mathbf{x}_j)$

for SVM; margin max hyperplane  $\mathbf{x}_i^T \mathbf{x}_j$

for log reg; min logistic loss  $\mathbf{x}_i^T \mathbf{x}_j$

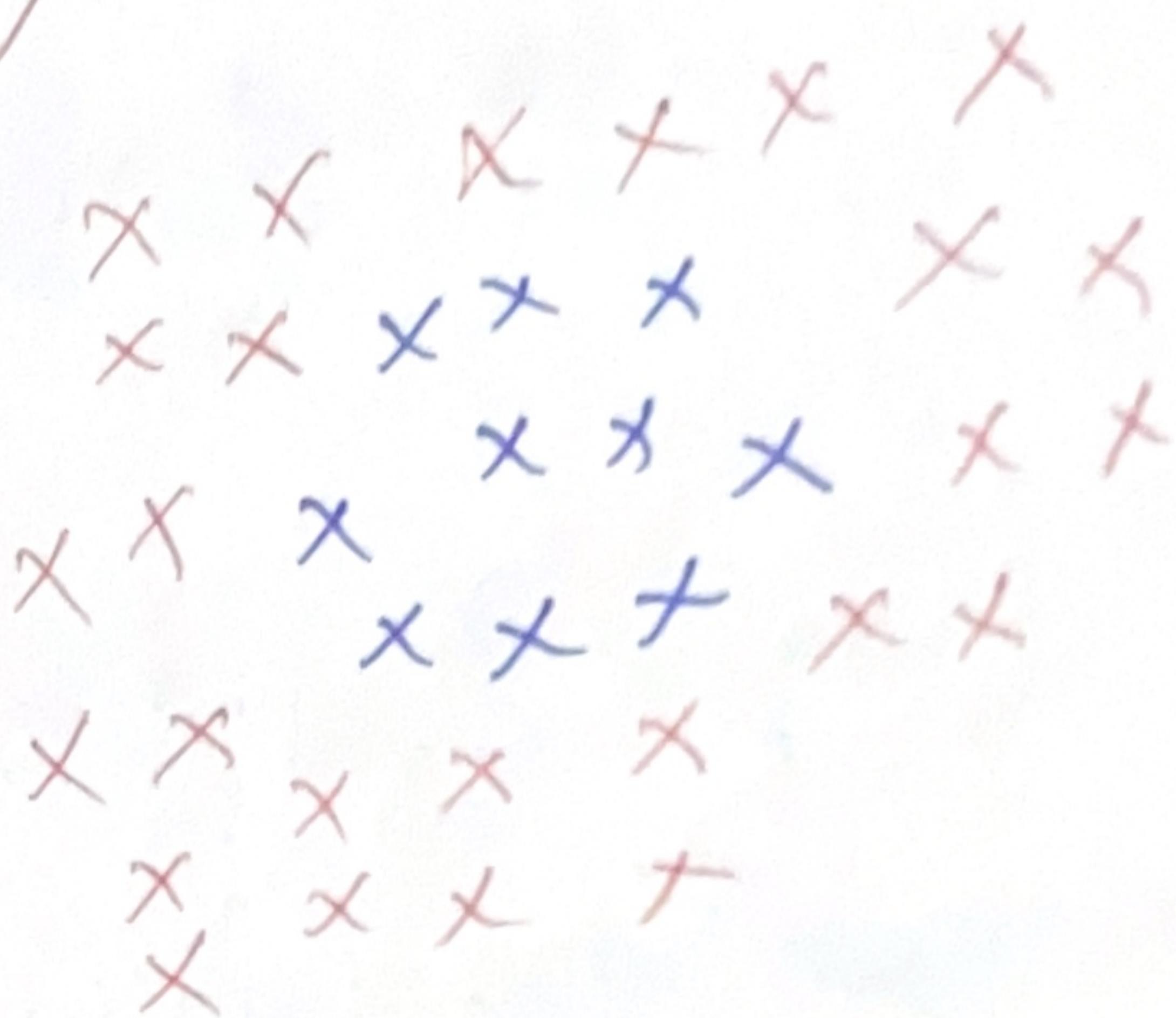
Model changing idea; Kernelization → 1990's

## Polynomial Kernel

→ Kernelization

$$K(x_1, x_2) = (x_1^T x_2 + c)^d$$

eg.  $K(x_1, x_2) = \underbrace{(1 + x_1^T x_2)^2}_{\text{Quadratic Kernel}}$



$$(f_1, f_2) \rightarrow (f_1^2, f_2^2)$$

$$K(x_1, x_2) = (1 + x_1^T x_2)^2$$

$$= (1 + x_{11}x_{21} + x_{12}x_{21} + x_{11}x_{22})^2 \quad x_1 = \langle x_{11}, x_{12} \rangle$$

$$x_2 = \langle x_{21}, x_{22} \rangle$$

logistic

$$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 + x_{22}^2 + 2x_{11}x_{21} + 2x_{12}x_{22} + 2$$

det =  $[1, x_{11}^2, x_{12}^2, \sqrt{2}x_{11}, \sqrt{2}x_{12} + \sqrt{2}x_{11}x_{12}, x_{11}x_{21}x_{12}x_{22}] : x_1^1$

$$= [1, x_{21}^2, x_{22}^2, \sqrt{2}x_{21}, \sqrt{2}x_{22}, \sqrt{2}x_{21}x_{22}] : x_2^1$$

$$= (x_1^1)^T (x_2^1)$$

$$\begin{matrix} x_1 \\ \downarrow \\ x_{11} \end{matrix} \quad \begin{matrix} x_2 \\ \downarrow \\ x_{21} \end{matrix}$$

$$\begin{matrix} x_1^T x_2 \\ \downarrow \\ x_1^T x_{21} \end{matrix}$$

Kernallization:  $d \xrightarrow[\text{internally}]{\text{feature trans}} d'$   $d' > d$

Mercuis Thm

Kernaltuck



$$(2d) \quad d \rightarrow d' (6d)$$

$d' > d$

not-eu-sep

$$K(x_1, x_2)$$

$\xrightarrow[2D]{(x_1, x_2)} \xrightarrow[6D]{(x'_1, x'_2)}$

Radial Basis Function (RBF) Kernel

SVM: most popular / general purpose: RBF

$$(x_1, x_2) \quad K_{RBF}(x_1, x_2) = \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma^2}\right)$$

c: hyperpara-

$$\begin{array}{ccc} & \overbrace{\hspace{1cm}}^{d_{12}} & \\ x_1 & & x_2 \end{array} \quad \|x_1 - x_2\|^2 = d_{12}^2$$

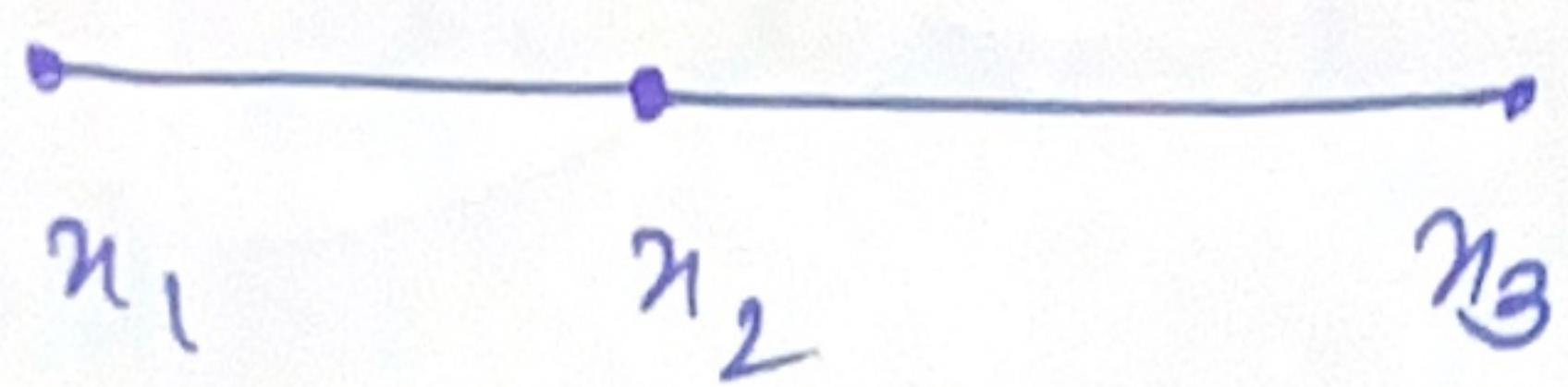
Soft margin Kernel SVM

(RBF)

c: hypuparameter

$$K(x_1, x_2) = \exp\left(-\frac{d_{12}^2}{2\sigma^2}\right) \quad d_{12} = \|x_1 - x_2\|$$

①  $d_{12} \uparrow: K(x_1, x_2) \downarrow$



$$K(x_1, x_2) > K(x_1, x_3)$$

$d \uparrow, d^2 \uparrow$

$e^{d^2} \uparrow$

$\frac{1}{e^{d^2}} \downarrow$

②  $\sigma$

$$\sigma = 1$$

$$\sigma = 0.1$$

$$\sigma = 10$$

$$\textcircled{1} \quad d=0$$

$$k=1$$

$\textcircled{2} \quad d \uparrow k \downarrow$

$$d(x_1, x_2) = 4 \\ K=0$$

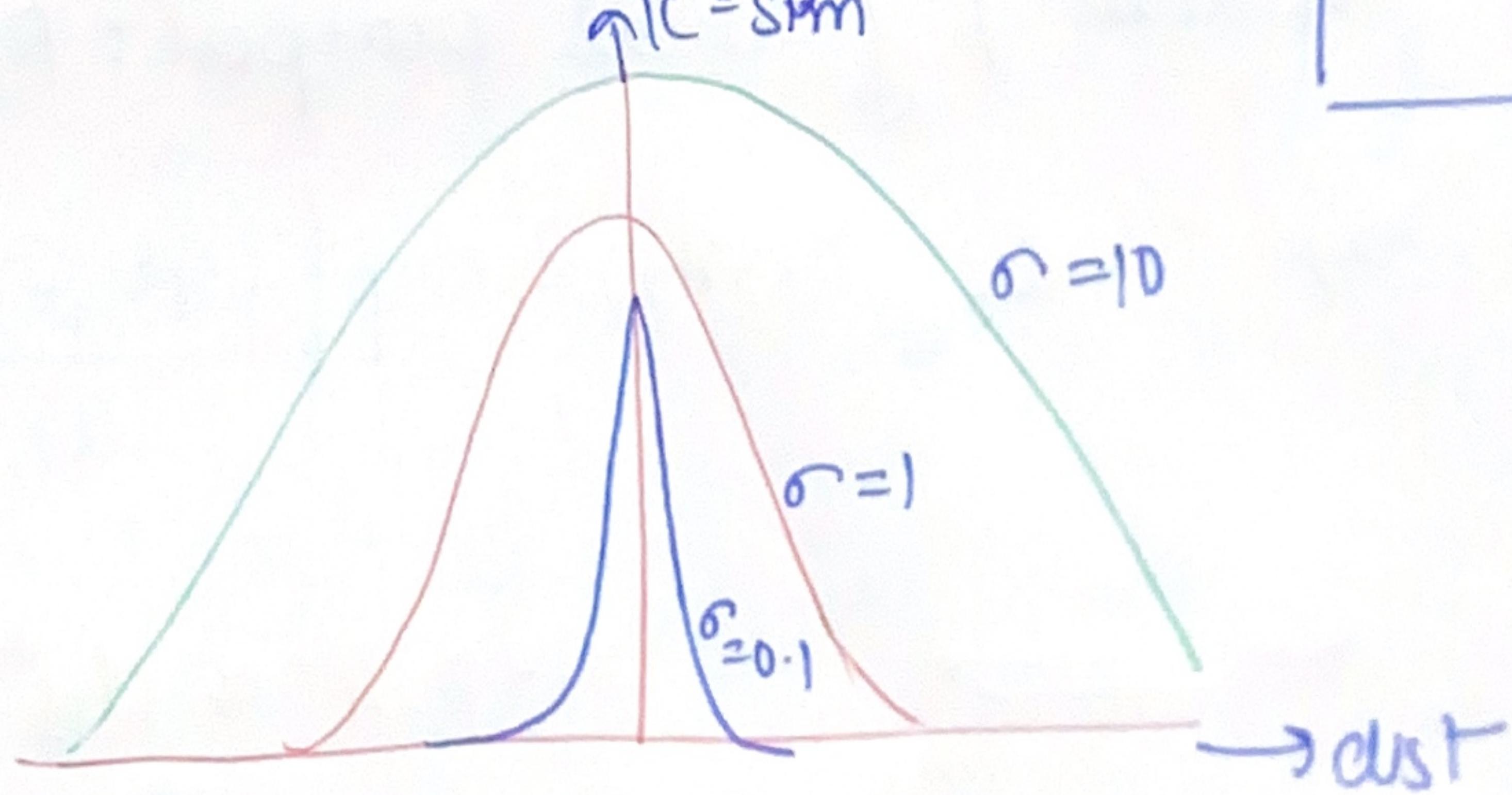
Kernel  $\rightarrow$  Sim func

distance  $\rightarrow$  dissimil func

$$K = \text{Sim}$$

$$d(x_1, x_2) = 0.5$$

$$\boxed{K=0 \rightarrow \text{Sim}=0 \\ \Rightarrow \text{dist}=\infty}$$



$\sigma \uparrow;$

$\begin{bmatrix} \text{dissim} \Leftrightarrow \text{dist} \\ \text{Sim} \Leftrightarrow K \end{bmatrix}$

## K-NN & RBF Kernel

RBF-SVM  
v K-NN

$\sigma \uparrow \Rightarrow K \uparrow$  in KNN  
(in RBF)

don't know the best Kernel



Simply use RBF SVM

Soft margin  $\leftarrow : C$        $\sigma : \text{RBF}$

$C, \sigma \leftarrow \text{grid/ Random Search}$

## Domain Specific kernels :

Polynomial;  $\underbrace{\text{RBF}}_{\hookrightarrow \text{KNN}} \rightarrow$  general purpose kernel

Kernel Trick  $\curvearrowright F^+ \rightarrow$  Domain specific  
aut

given problem

$\hookrightarrow$  string kernels  $\rightarrow$  Text classfr

FT: hard part of ML

$\hookrightarrow$  partially replaced by finding right kernel

# Train and Run Time Complexity of SVM

Train → SGD

↓  
specialized alg (dual) → Sequential  
minimal  
optimization  
(SMO)

LibSVM: best library for training SVMs ↑

Sklearn

Training Time:  $\sim O(n^2)$  for kernel SVMs

$\nu(2007)$ :  $O(nd^2)$  if  $d < n$

if  $n$  is large  $\rightarrow O(n^2) \uparrow \uparrow$

↑ Typically do not use SVM when  $n$  is large

↓  
internet

Run Time  $f(x_q) = \sum_{i=1}^n \alpha_i y_i k(x_i, x_q) + b$

$\alpha_i = 0$  for non SVs

# SVs =  $k$

$O(kd)$  if # SVs are small  
 $O(kd)$  is small

## nu-SVM

C-SVM  $\rightarrow$  original formulation  $[C \geq 0]$

parameter

alternate formulation of SVM

$\text{nu-SVM}$   $0 \leq \text{nu} \leq 1$

hyperparameter  $\text{nu} \geq \text{fraction of errors}$

$\text{nu} \leq \text{fraction of SVs}$

$\text{SVM}$  Dtrain  
 10% errors  
 $\text{nu} = 0.1$   
 1% errors  
 $\text{nu} = 0.01$

$\text{nu} = 0.01 \Rightarrow \% \text{ age of errors} \leq 1\%$   
 $\#\text{SVs} \geq 1\% \text{ of } n$

Run Time complex : fewer SVs

But

$n = 0.01 \Rightarrow \text{errors} \leq 1\%$   
 but  $\#\text{SVs} \geq 1\% \text{ of } n$

$n = 100,000$   
 $\#\text{SVs} = k \geq 1000$

Support Vector regression (SVR)  $y_i \in \mathbb{R}$

SVM classfn SVC  $\rightarrow y_i \in \{+1, -1\}$

Math

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

$\hat{y}_i$

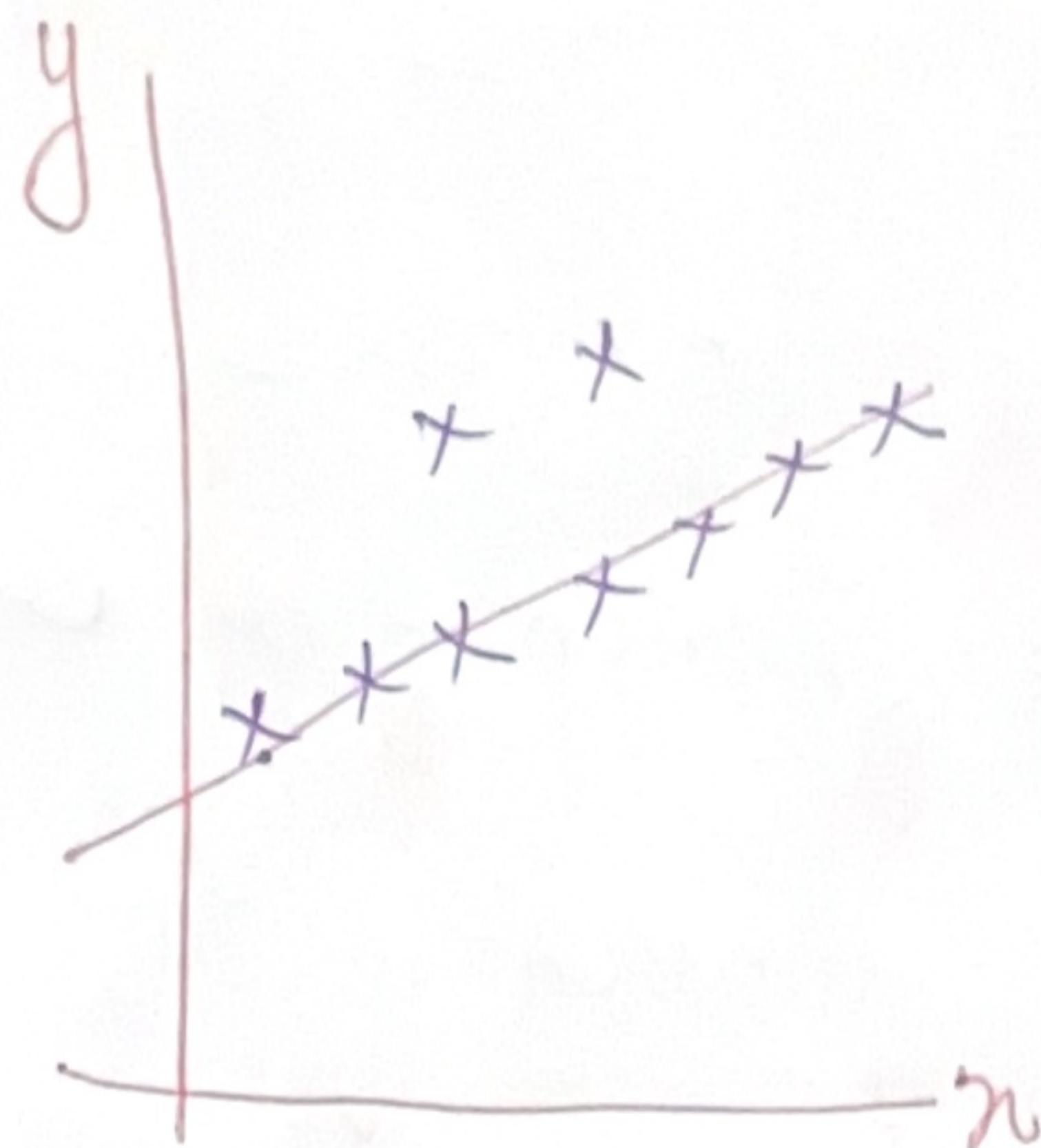
hypurp ] linear form

sat  $y_i - (w^T x_i + b) \leq \epsilon \rightarrow$  SVR

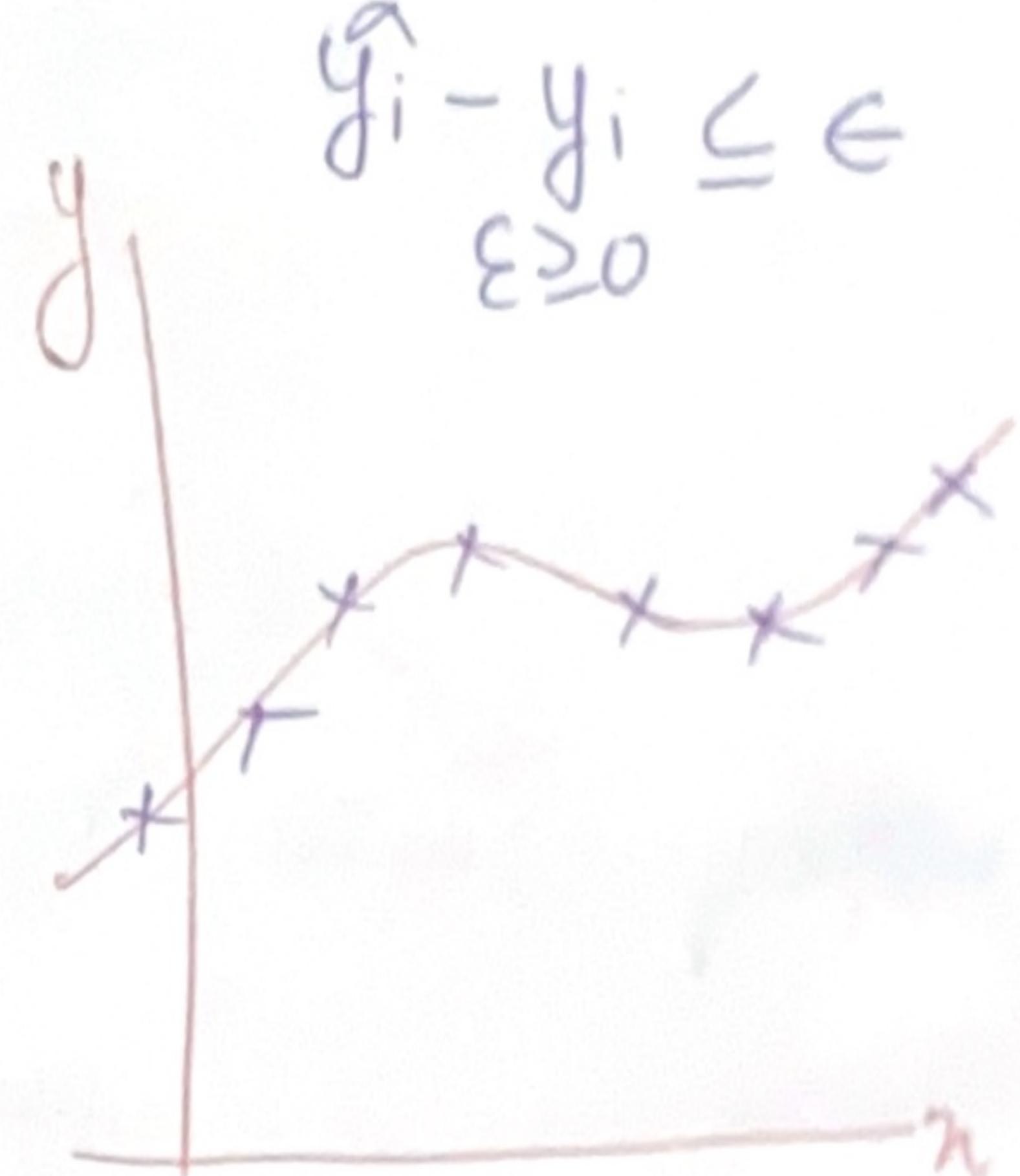
$$(w^T x_i + b) - y_i \leq \epsilon \quad \forall i$$

$\epsilon \geq 0$

$$\left. \begin{aligned} f(x_i) &= w^T x_i + b \\ &= \hat{y}_i \\ y_i - \hat{y}_i &\leq \epsilon \end{aligned} \right\}$$



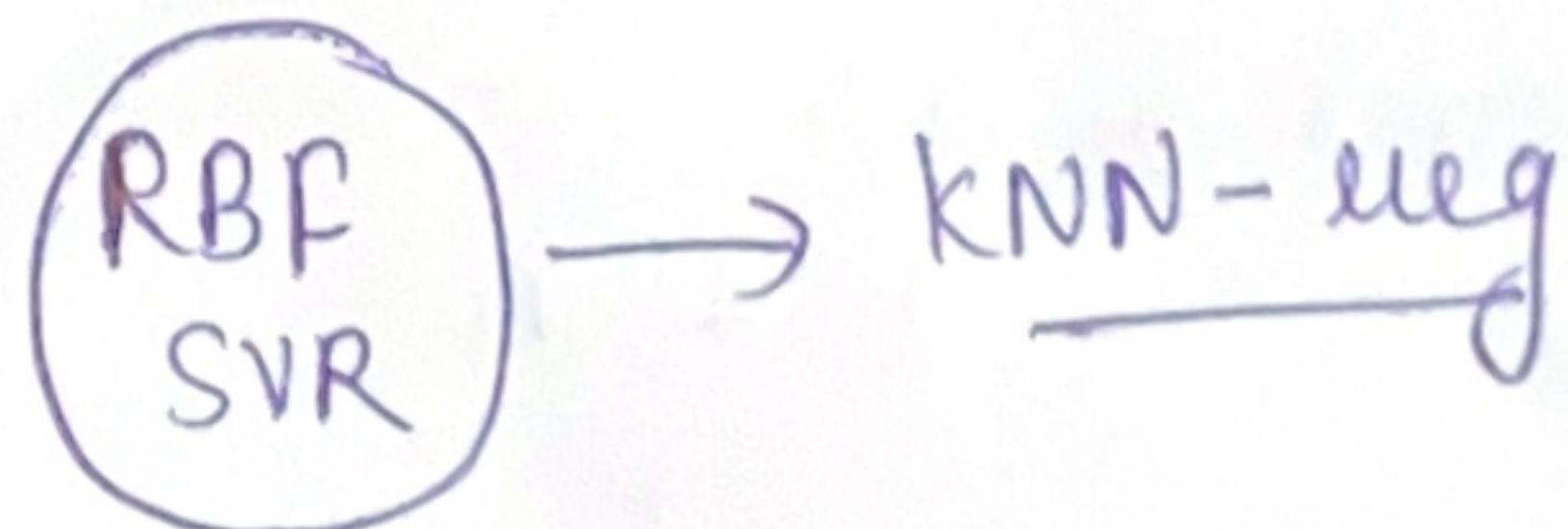
Linear SVR



Kernel SVR

$\epsilon \downarrow \Rightarrow$  errors are low on training data  
 $\Rightarrow$  overfitting  $\uparrow$

$\epsilon \uparrow \Rightarrow$  errors on  $\mathcal{D}_{\text{train}} \uparrow \Rightarrow$  underfitting  $\uparrow$



### Cases (SVM)

$\rightarrow$  feature Engg & FT: kernel design, finding the right kernel

$\rightarrow$  decision Surface for SVM: hyperplane

Kernel SVM: (a)  $x_i \rightarrow$  non-linear  
 $\downarrow$  kernel trick      Surface  
 $d' \gg d$       (d')       $x_i' \rightarrow$  linear Surface

- Similarity / dist func
  - ↪ Kernel:  $K(x_i, x_j)$
- Interpretability & feature importance  $\rightarrow$  Kernel SVM
  - ↪ FI for various features
  - ↪ forward feature selection
- Outliers  $\rightarrow$  v. little impact
  - ↪ SVs that matter
  - ↪ RBF with small  $\sigma$ ;  $\rightarrow$  KNN with small  $k$
- Bias-var:  $C \uparrow \rightarrow$  overfit  $\Rightarrow$  high-var
  $C \downarrow \rightarrow$  underfit  $\Rightarrow$  high bias
   
 RBF-SVM:  $C, \sigma \downarrow \rightarrow$  RBF
- Large  $d \rightarrow$  v. good for SVM
  - $\textcircled{d} \leftarrow \textcircled{d}^{\text{II}}$
  - ↪ Kernel (RBF)
- Best case: (tight kernel)

Worst case  $n$  is large  $\rightarrow$  Train time is high

Logistic regression  $K$  is large  $\rightarrow$  low latency is not possible  
 $\downarrow$  # SVs