

## **Last time**

Perceptron

Support Vector Machines

Primal Problem – Margin Maximization

Dual Problem

## **Today – part 1**

- - Kernel Version

- What if there are outliers in the problem?

$$\begin{aligned} \min_w \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t. } & y_i (w^T x_i) y_i \geq 1 \end{aligned}$$

### DUAL PROBLEM

can be  
KERNELIZED!

Solving in  
n instead of d.

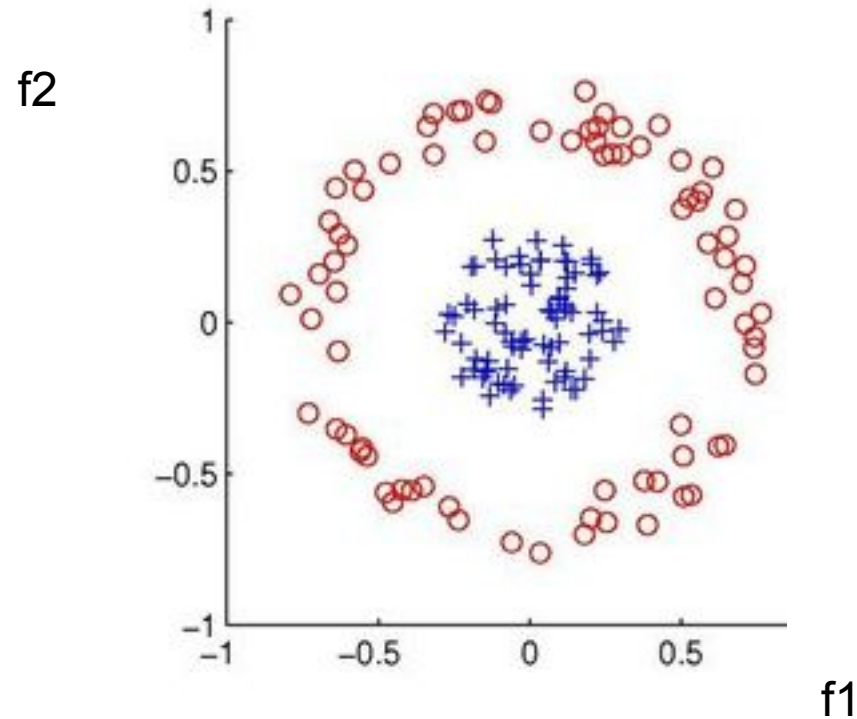
$$\max \quad \alpha \geq 0$$

easy  
constraints

$$\alpha^T \mathbf{1} - \frac{1}{2} \alpha^T y^T \underbrace{X^T X}_{n \times n} y \alpha$$

n x n.

ISSUE → Features could be non-linearly related



In general,

$$(f_1 - a)^2 + (f_2 - b)^2 = r^2$$

Every datapoint  $x$   
roughly satisfies

$$\underline{f_1^2} + \underline{a^2} - \underline{2f_1 a} + \underline{f_2^2} + \underline{b^2} - \underline{2f_2 b} - \underline{r^2} = 0$$

$$\begin{bmatrix} f_1 & f_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & f_1 & f_2 & f_1 f_2 & f_1^2 & f_2^2 \end{bmatrix}$$

$$x \in \mathbb{R}^2 \rightarrow \underline{\phi(x)} \in \mathbb{R}^6$$

Datapoints  $\phi(x)$

satisfies

$$\begin{bmatrix} 1 & f_1 & f_2 & f_1 f_2 & f_1^2 & f_2^2 \end{bmatrix}$$

$\phi(x)$

$u \in \mathbb{R}^6$

$$\boxed{\phi(x)^T u = 0}$$

$$\begin{bmatrix} \underline{a^2 + b^2 - r^2} \\ -2a \\ -2b \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Idea: Transform features from

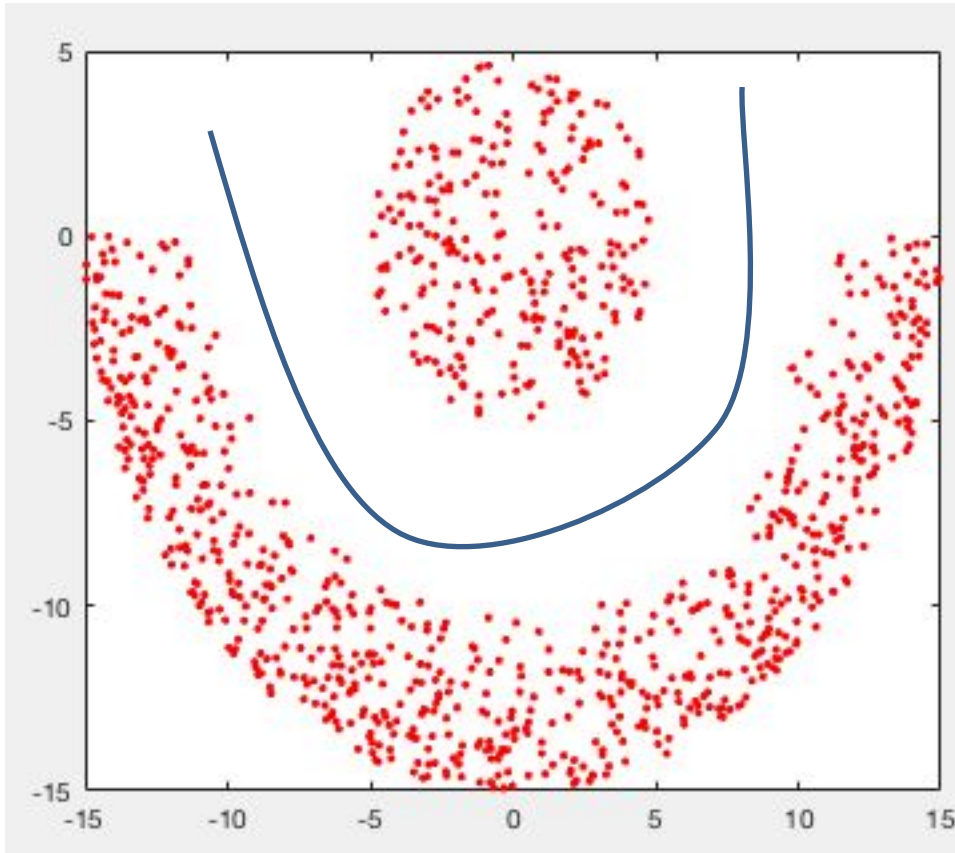
low dimension  $\mathbb{R}^d$  to

high dimension  $\mathbb{R}^D$

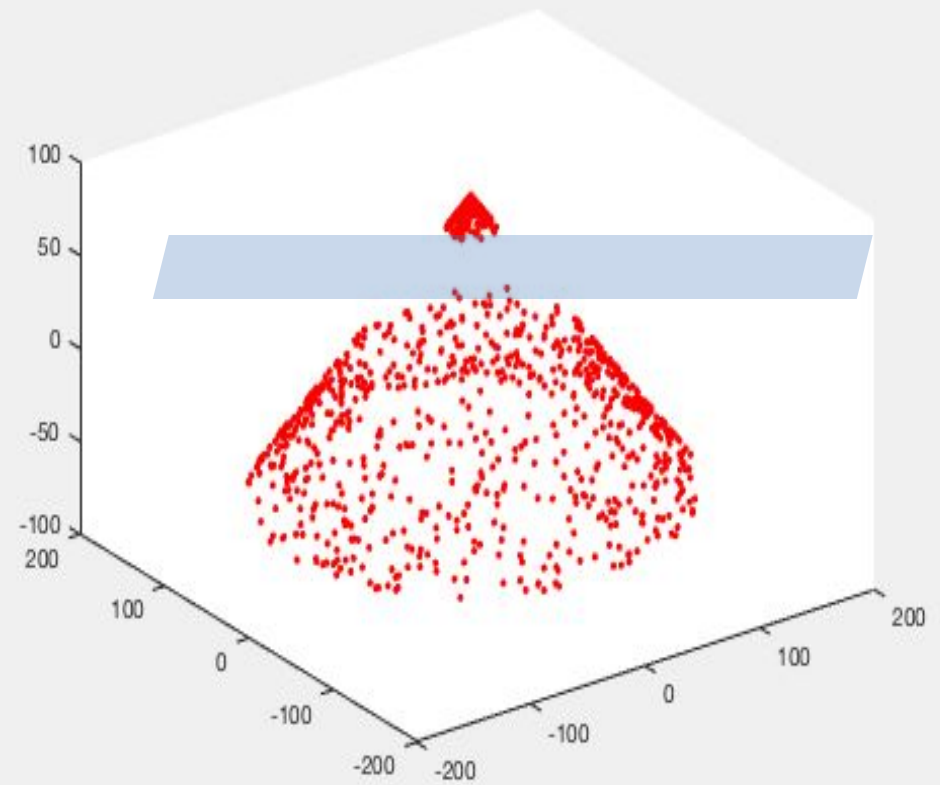
$$x \rightarrow \phi(x)$$

$\mathbb{R}^d \quad \mathbb{R}^D$

ORIGINAL DATA



3D REPRESENTATION



Quadratic Map  $\phi$   
to 6 Dimensions

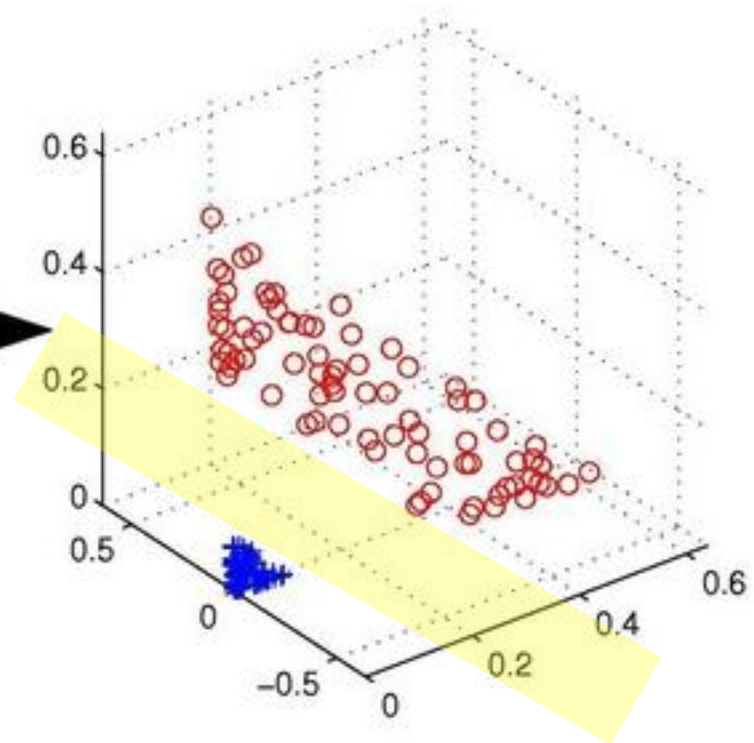
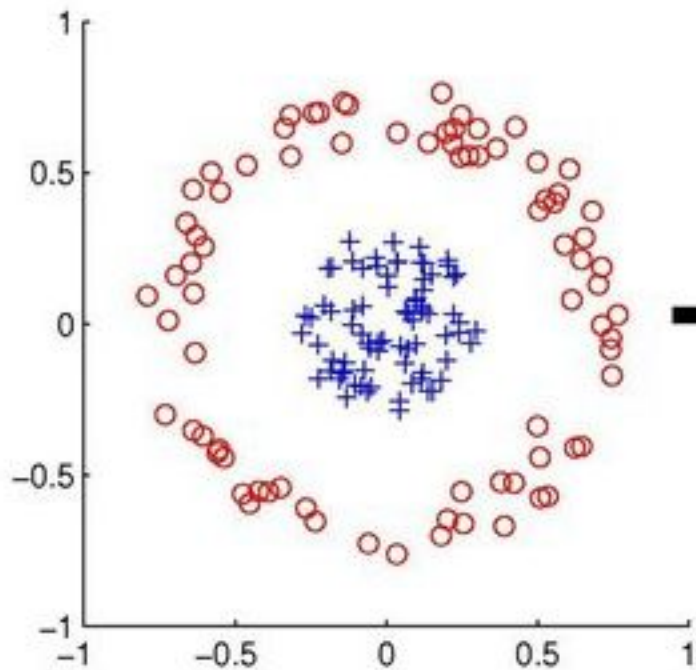
PCA

$w_1, w_2, w_3$   
Top 3 principal  
components

Reconstruct  
each point  $x$  as  
 $[\phi(x)'w_1 \ \phi(x)'w_2 \ \phi(x)'w_3]$

$$\{x_1, x_2, \dots, x_n\}$$

$$\{\Phi(x_1), \Phi(x_2), \dots, \Phi(x_n)\}$$



$$\forall i, x_i \in \mathbb{R}^2$$

$$\forall i, \phi(x_i) \in \mathbb{R}^3$$

$$x = [f_1 \quad f_2 \quad f_3 \quad f_4]$$

$\Downarrow$  Cubic relation

$$\phi(x) = \left[ 1 \quad f_1 \quad f_2 \quad f_3 \quad f_4 \right]$$

$\underbrace{\hspace{10em}}$   
 $1 + 4 \quad + \quad 4c^2 + 4c^3$

$$x \in d$$

$$\phi(x) \in d(d^p)$$

ISSUE :

$\phi(x) \in \mathbb{R}^D \rightarrow$  might be too large



# EXAMPLE

$$x = [f_1 \quad f_2]$$

$$x' = [g_1 \quad g_2]$$

$$\begin{aligned} \boxed{\left( \underline{x x' + 1} \right)^2} &= \left( [f_1 \quad f_2] \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} + 1 \right)^2 \\ &= \left( f_1 g_1 + f_2 g_2 + 1 \right)^2 \\ &= \underline{f_1^2 g_1^2} + \underline{f_2^2 g_2^2} + 1 + 2 f_1 g_1 f_2 g_2 \\ &\quad + \underline{2 f_1 g_1} + 2 f_2 g_2 \end{aligned}$$

$$\left( \underline{\underline{\tau x' + 1}} \right)^2$$

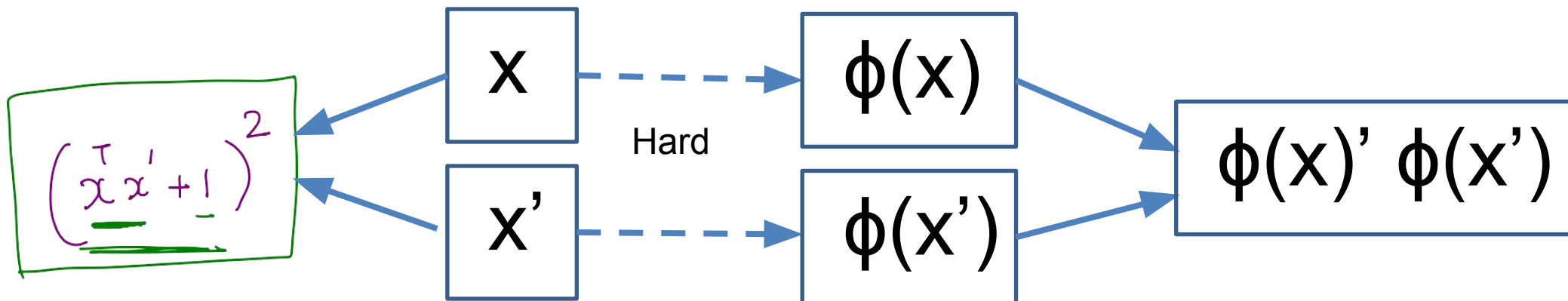
$$= \left[ \underline{f_1^2}, \underline{f_2^2} \right]$$

$$\uparrow \phi(x)^\tau \phi(x')$$

$$\sqrt{2} f_1 f_2$$

$$\underline{\sqrt{2} f_1}$$

$$\sqrt{2} f_2 \begin{bmatrix} \underline{g_1^2} \\ \underline{g_2^2} \\ 1 \\ \sqrt{2} g_1 g_2 \\ \sqrt{2} g_1 \\ \sqrt{2} g_2 \end{bmatrix}$$



## So far

- > To capture non-linear relationships, one can “create” non-linear functions of features.
- > But the number of features to create grows exponentially with the degree of non-linearity  $p$  that we wish to capture ( $d^p$ )
- > For  $d = 2$  and  $p = 2$ , it appears there is a trick to get around this.
- > Is this trick general enough to be useful the general case as well? (i.e., for any  $d$  and any  $p$ ?)

## MORE EXAMPLES

Polynomial map

$$k(x, x') = \left( \underline{x^T x' + 1} \right)^p$$

for some  $p \geq 1$

→ Can be shown to be a "valid" function.

i.e.,  $\exists \phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$  such that

$$k(x, x') = \phi(x)^T \phi(x')$$

EXERCISE.

Compute  $\phi$   
for  $p=3$ ,  
 $p=4$ .

## MORE EXAMPLES

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

for some  $\sigma > 0$



RADIAL BASIS FUNCTION

→ Can also shown to be a "valid" map

→ Interestingly,  $\phi$  in this case maps  $x$  to an "infinite" dimensional space

[ Technicalities aside  
think of  $\phi(x)$  as a  
function and dot product as  
integrals ]

## KERNEL FUNCTION

Any function  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  which is a "valid" map is a kernel function

$$K(x, x') = (x^T x' + 1)^p \rightarrow \text{Polynomial kernel}$$

$$K(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \rightarrow \text{Gaussian kernel / Radial basis kernel}$$

Question:

Given a function

$\mathbb{R} \cdot \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ , how can we say it's a

valid kernel?

METHOD 1.

Explicitly construct a  $\phi$  map

[might be hard sometimes]

## METHOD 2: MERCER'S THEOREM

A function  $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a valid kernel  
if and only if

(a)  $k$  is symmetric.

(b) For any dataset  $\{x_1, \dots, x_n\}$ , the

matrix  $K \in \mathbb{R}^{n \times n}$

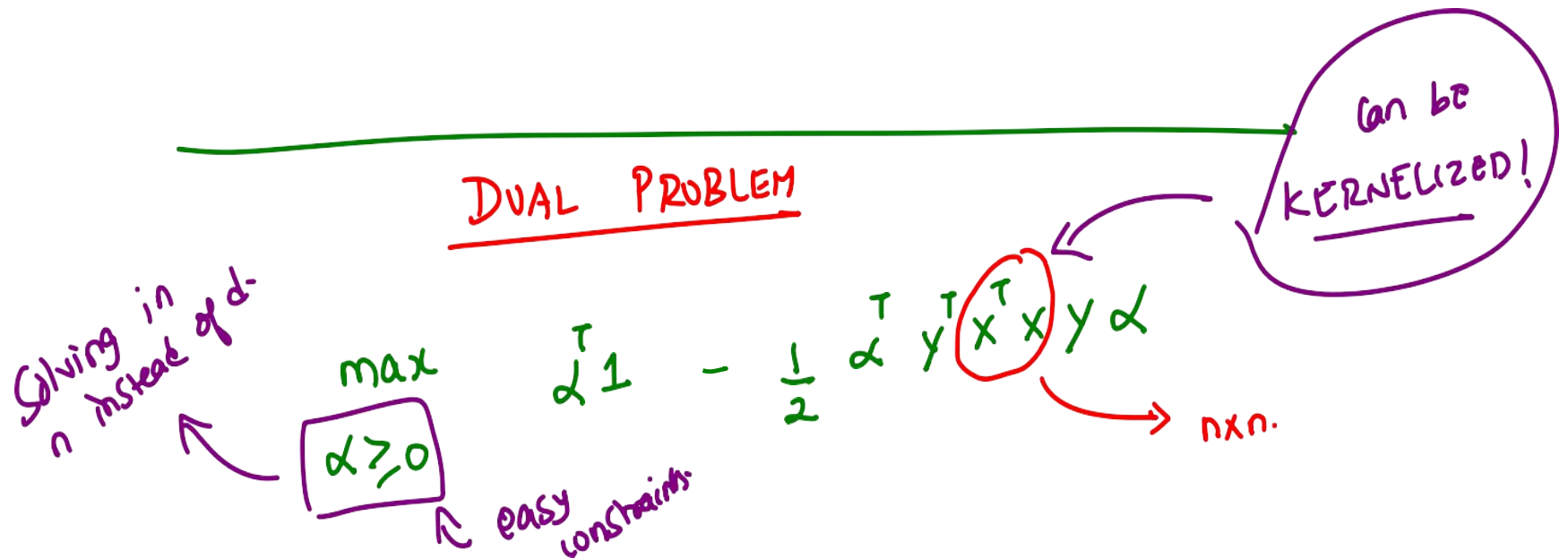
$$K_{ij} = k(x_i, x_j)$$

is

POSITIVE  
SEMI  
DEFINITE

eigenvalues  
of  
 $K$  are all  
non-negative.





Instead of  $X^T X$ , use  $K$  for the chosen Kernel.

- Cannot recover  $w = \sum \alpha_i \phi(x_i)$  – *but that's okay. Why?*
- Solve the dual problem. Obtain alphas.
- For a test point, to predict, use the following:

$$w^T \Phi(x_{\text{test}}) = \left( \sum \alpha_i \phi(x_i) \right)^T \Phi(x_{\text{test}}) = \sum \alpha_i K(x_i, x_{\text{test}})$$

Idea (to deal with outliers):

Fix any  $w$ .  $w$  classifies some points

correct and some incorrectly. Let the

incorrect points pay "bribe" to get to the

correct side.

Modified formulation

$$\min_{\omega, \xi} \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i$$

$C \geq 0$  [hyper parameter]

$$\rightarrow (\omega^T x_i) y_i + \underline{\xi_i} \geq 1 \quad \leftarrow \forall i$$

$$\rightarrow \underline{\xi_i} \geq 0 \quad \leftarrow \forall i$$

if  $C = 0 \Rightarrow$  Bribes don't cost  $\Rightarrow$   $\omega = 0$  is solution

$C \rightarrow \infty \Rightarrow$  Bribes are too costly  $\Rightarrow$  Linear separable cost.

$$L(w, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + c \underbrace{\left( \sum_{i=1}^n \xi_i \right)}_{\uparrow} + \underbrace{\sum_{i=1}^n \alpha_i \left( 1 - \underbrace{(w^T x_i) y_i}_{-\xi_i} \right)}_{\text{green underline}} + \sum_{i=1}^n \beta_i \underbrace{(-\xi_i)}_{\uparrow}$$

Dual:

$$\max_{\substack{\alpha \geq 0 \\ \beta \geq 0}} \min_w L(w, \xi, \alpha, \beta)$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w^* = \sum_{i=1}^n \alpha_i x_i y_i$$

$$w^* = x y \alpha$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \boxed{C - \alpha_i - \beta_i = 0}$$

$$\boxed{\alpha_i + \beta_i = C} \quad \forall i$$

Substitute  $w^* = x\gamma\alpha$  in the original objective

$$\frac{1}{2} (x\gamma\alpha)^T (x\gamma\alpha) + \sum_{i=1}^n \underbrace{(C - \alpha_i - \beta_i)}_{=0} \xi_i + \alpha^T \mathbf{1} - (x\gamma\alpha)^T (x\gamma\alpha)$$

SOFT-MARGIN  
SUPPORT  
VECTOR  
MACHINE

$$\begin{array}{l} \max \\ \alpha \geq 0 \\ \beta \geq 0 \\ \hline \alpha + \beta = C \end{array} \quad \frac{\alpha^T 1 - \frac{1}{2} (xy\alpha)^T (xy\alpha)}{}$$

$\equiv$

$$\begin{array}{l} \max \\ 0 \leq \alpha \leq C \\ \hline \end{array} \quad \frac{\alpha^T 1 - \frac{1}{2} \alpha^T y^T (x^T x) y \alpha}{}$$

Box  
CONSTRAINT.

# Summary

HARD-MARGIN  
SVM

PRIMAL

$$\min_w \frac{1}{2} \|w\|^2$$

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

$1 - w^T x_i y_i \leq 0$

DUAL

$$\max_{\alpha \geq 0} \alpha^T 1 - \alpha^T y^T \underline{x^T x} y$$

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

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

SOFT-MARGIN  
SVM

PRIMAL ✓

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

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

$\alpha \rightarrow \xi_i \geq 0 \quad \forall i = 1, \dots, n$

DUAL ✓

$$\max_{\alpha, \beta} \alpha^T 1 - \alpha^T y^T \underline{x^T x} y$$

$\alpha + \beta = C$   
 $\alpha \geq 0$   
 $\beta \geq 0$

$$0 \leq \alpha \leq C$$

- Let  $(\underline{w}^*, \underline{\xi}^*)$  be the primal optimal solution
- Let  $(\underline{\alpha}^*, \underline{\beta}^*)$  be the dual optimal solution

### COMPLEMENTARY SLACKNESS

$$\underline{\alpha}_i^* \left( 1 - \underline{w}^{*T} x_i - \underline{\xi}_i^* \right) = 0 \quad \forall i$$

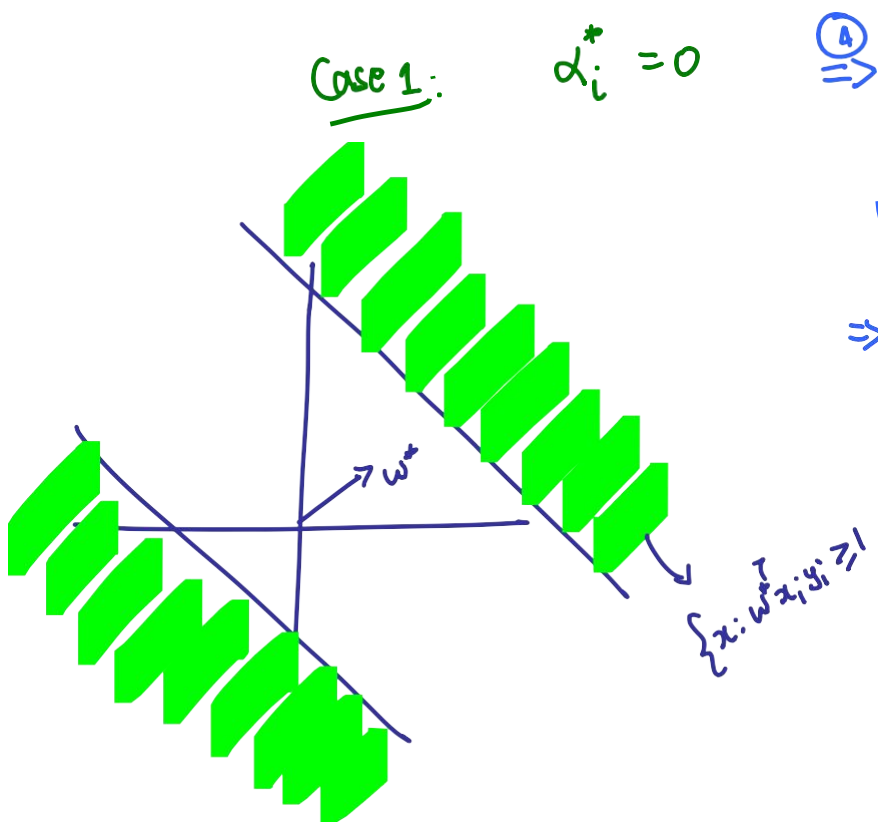
$$\underline{\beta}_i^* \underline{\xi}_i^* = 0 \quad \forall i$$

$$\underline{\alpha}_i^* + \underline{\beta}_i^* = c \quad \forall i$$

↳ (A)

Various cases possible





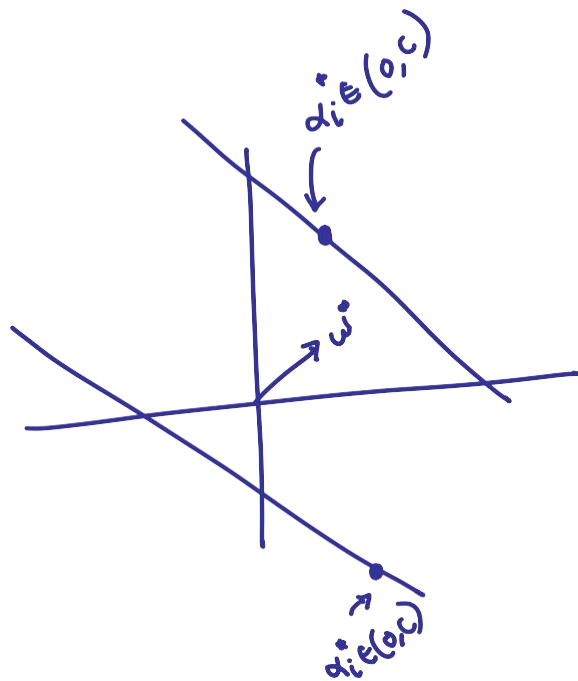
$$\stackrel{\textcircled{A}}{\Rightarrow} p_i^* = C \stackrel{\boxed{CS}}{\Rightarrow} \underline{\xi_i^*} = 0$$

$$1 - (w^*{}^T x_i) y_i - \underline{\xi_i^*} \leq 0 \quad [\text{Primal feasibility}]$$

$$\Rightarrow 1 - (w^*{}^T x_i) y_i \leq 0$$

$$\Rightarrow w^*{}^T x_i y_i \geq 1$$

$\Rightarrow w^*$  classifies  $(x_i, y_i)$  correctly.



Case 2:

$$0 < \alpha_i^* < C \quad \textcircled{A} \Rightarrow$$

$$0 < \beta_i^* < C \quad \stackrel{CS}{\Rightarrow} \underline{\xi_i^*} = 0$$

$$\Downarrow \boxed{CS}$$

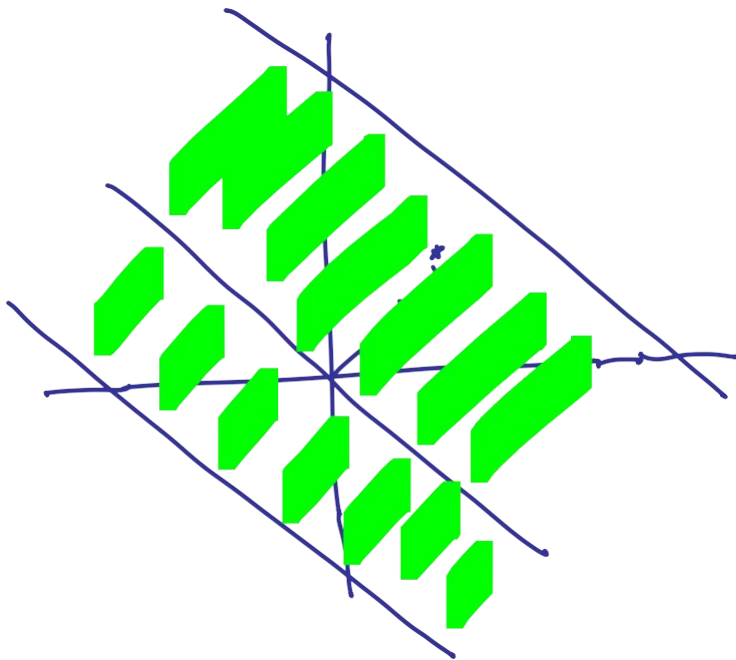
$$1 - (w^T x_i) y_i - \xi_i^* = 0$$

$$\Downarrow$$

$$(w^T x_i) y_i = 1$$

$$\Rightarrow$$

$(x_i, y_i)$  lies on the  
Supporting hyperplane.



Case 3:  $\alpha_i^* = C \Rightarrow p_i^* = 0 \Rightarrow \xi_i^* \geq 0$

$\Downarrow$  CS

$$1 - w^T x_i y_i - \xi_i^* = 0$$

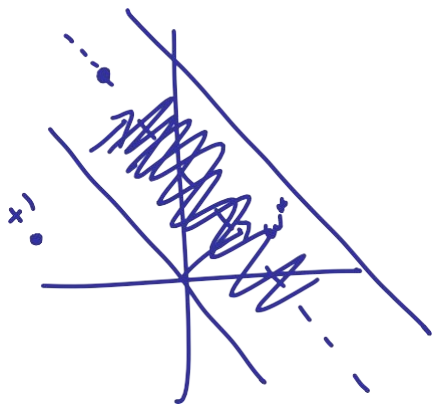
$$\xi_i^* = 1 - w^T x_i y_i \geq 0$$

$$\Rightarrow \boxed{w^T x_i y_i \leq 1}$$

Lets see this from P.O.V of data

CASE 1

$$\boxed{w^* x_i y_i < 1}$$



$$1 - w^* x_i y_i - \xi_i^* \leq 0$$

$$w^* x_i y_i \geq 1 - \xi_i^*$$

$$\boxed{\xi_i^* \geq 1 - w^* x_i y_i}$$

$$\Rightarrow \xi_i^* > 0 \Rightarrow \beta_i^* = 0 \Rightarrow \underline{\alpha_i^* = C}$$

$$\begin{aligned} \underline{\alpha_i^*} (1 - w^* x_i y_i - \xi_i^*) &= 0 \\ \beta_i^* \xi_i^* &= 0 \end{aligned}$$

CASE 2:  $w^* x_i y_i = 1$

$$\xi_i^* \geq 1 - \underline{w^* x_i y_i}$$

$$\Rightarrow \xi_i^* \geq 0 \Rightarrow \alpha_i^* \in [0, c]$$

CASE 3  $w^* x_i y_i > 1$

$$1 - \underbrace{w^* x_i y_i} - \xi_i^* \leq 0 \quad [\text{Primal feasibility}]$$

$$\Rightarrow 1 - w^* x_i y_i - \xi_i^* < 0 \quad \Rightarrow \quad \boxed{\text{CS}} \quad \alpha_i^* = 0$$

