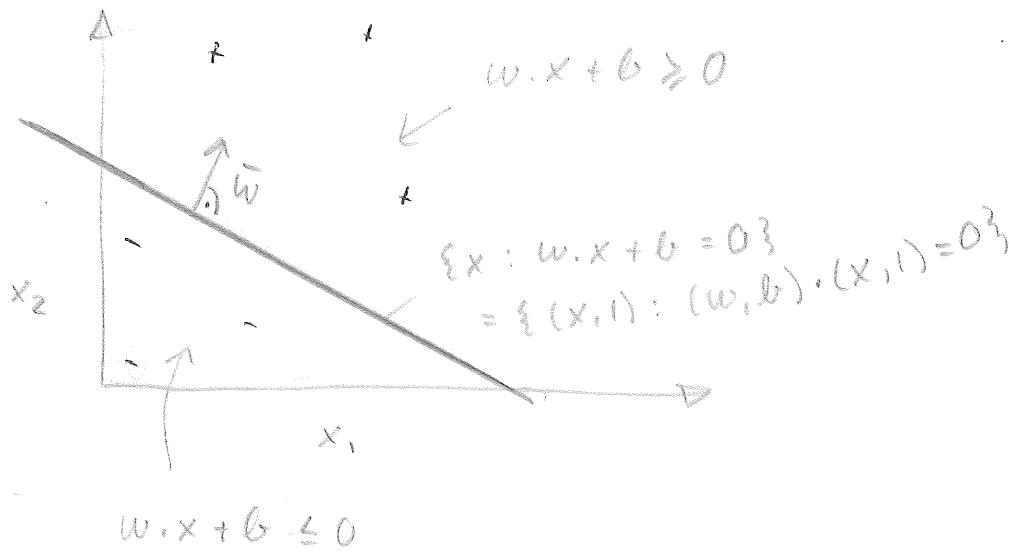


## SUPPORT VECTOR MACHINES

SEPARATE DATA INTO 2 CLASSES  
BASED ON A HYPER PLANE



IF  $b \geq 0$  THEN ORIGIN ON  $\geq$  SIDE

WHAT IS DISTANCE TO ORIGIN ?

INTERSECTION BETWEEN

LINE  $x\bar{w}$  & PLANE :

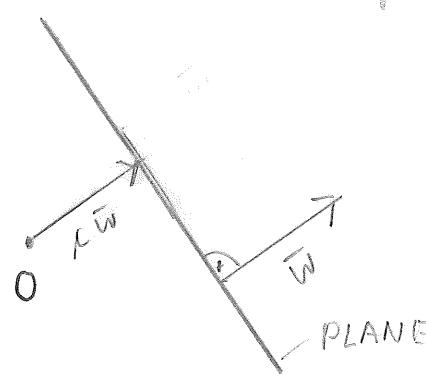
$$w \cdot (cw) + b = 0$$

$$\therefore c = \frac{-b}{\|w\|^2}$$

$$\|cw\| = |c|\|w\|$$

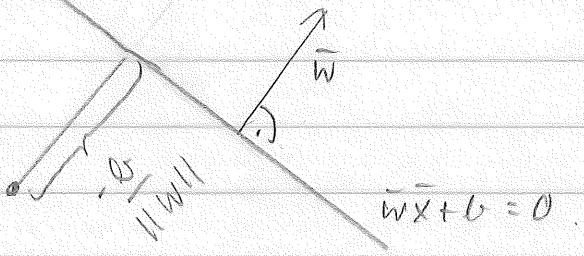
$$= \frac{|b|}{\|w\|}$$

DISTANCE



THE SIGNED OF ORIGIN TO  $\{x : w \cdot x + b = 0\}$   
is  $\frac{-b}{\|w\|}$ , ^ DISTANCE IS 0 IFF  $b = 0$

2



NOTE THAT  $\bar{w} \cdot \bar{x} + b = 0$

$$\Leftrightarrow q \bar{w} \cdot \bar{x} + q b = 0 \quad \text{FOR SCALAR } q \neq 0$$

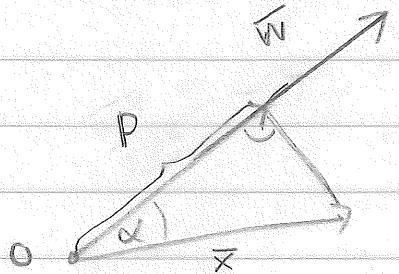
↑ SAME PLANE

$$\Leftrightarrow \frac{\bar{w}}{\|\bar{w}\|} \cdot \bar{x} + \frac{b}{\|\bar{w}\|} = 0$$

$\frac{\bar{w}}{\|\bar{w}\|}$   
UNIT

DOT PRODUCT:  $\bar{w} \cdot \bar{x} = \|w\| \|x\| \cos(\alpha)$

P IS PROJECTION OF  
X ONTO W

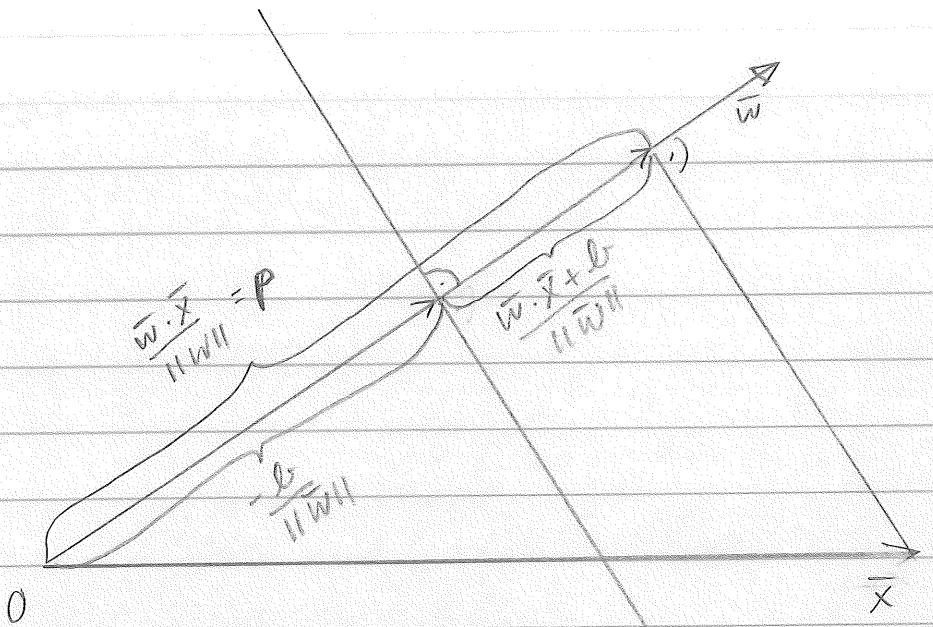


$$\cos(\alpha) = \frac{p}{\|x\|}$$

$$p = \|x\| \cos(\alpha)$$

$$p = \frac{\bar{w} \cdot \bar{x}}{\|\bar{w}\|}$$

2-3

PLANE  $(\bar{w}, b)$ 

$$\frac{|\bar{w} \cdot \bar{x} + b|}{\|\bar{w}\|}$$

DISTANCE OF  $\bar{x}$  TO PLANE

$$\frac{\bar{w} \cdot \bar{x} + b}{\|\bar{w}\|}$$

SIGNED DISTANCE OF PLANE TO  $\bar{x}$

WANT  $w \cdot x_i + b > 0$  IF  $y_i = +1$   
 $< 0$  IF  $-1$

---

WANT  $y_i (w \cdot x_i + b) > 0$

$\underbrace{w \cdot x_i}_{\gamma_i}$

FUNCTIONAL MARGIN OF EXAMPLE  $(\bar{x}_i, y_i)$   
 WRT  $(\bar{w}, b)$  IS  $\delta_i = y_i(\bar{w} \cdot \bar{x}_i + b)$

GEOMETRIC MARGIN

$$\text{acc } \frac{\delta_i}{\|\bar{w}\|}$$

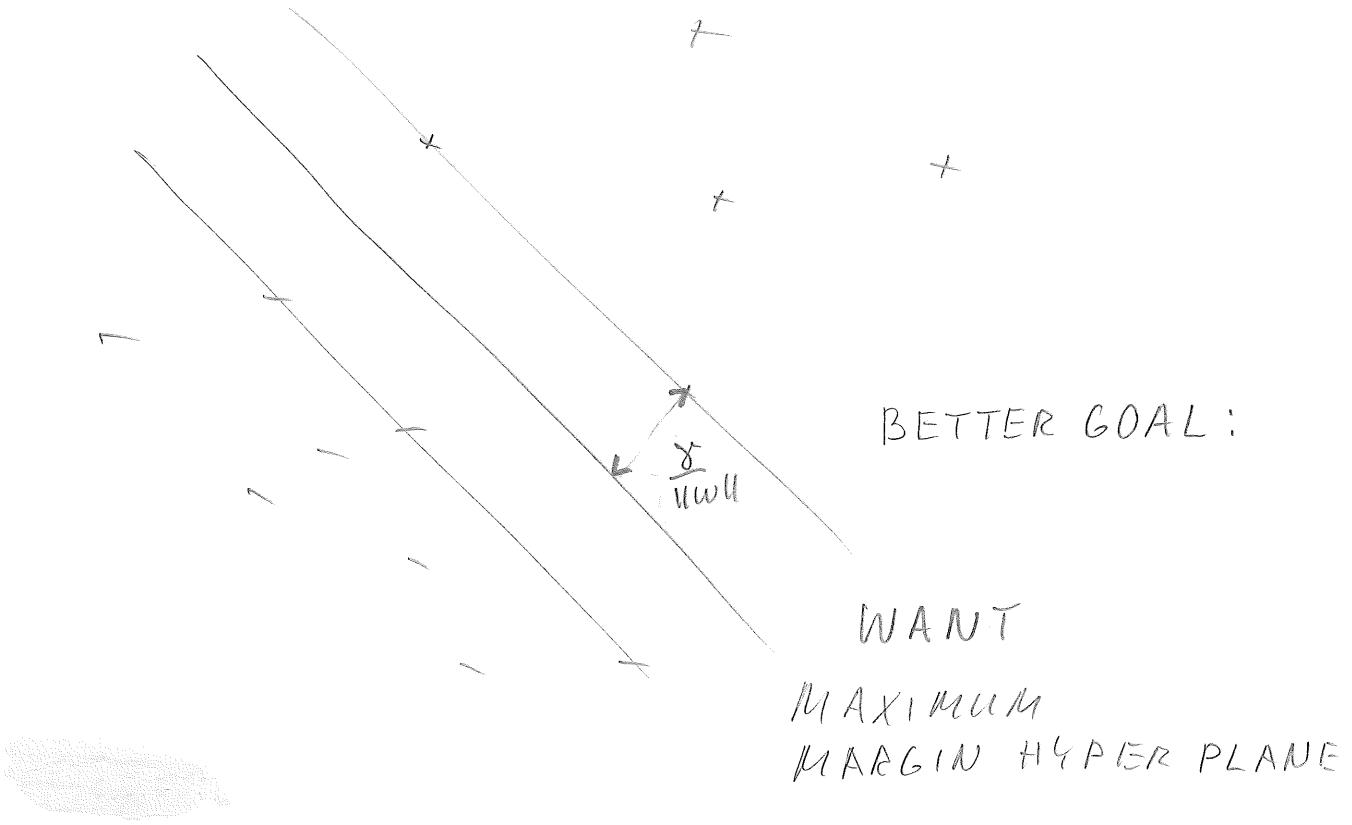
GEOMETRIC MARGINS ARE SCALED

DISTANCES

IF  $\|\bar{w}\| = 1 \Rightarrow$  NO NORMALIZATION NECESSARY.

WANT PLANE S.T. ALL  
 EXAMPLES HAVE POSITIVE MARGIN

4



PLANES  $(w, b)$  &  $(\lambda w, \lambda b)$  FOR  $\lambda > 0$  ARE THE SAME

GEOMETRIC MARGIN OF  $x_i$  INVARIANT TO SCALING

$$\frac{w}{\|w\|} x_i + \frac{b}{\|w\|}$$

HOW TO DEFINE MAX. MARGIN HYPERPLANE

$$\begin{aligned} & \text{MAX}_{w,b,\gamma} \quad \gamma \leftarrow \text{FUNCTIONAL} \\ & \text{SUBJECT TO} \quad y_i(w \cdot x_i + b) \geq \gamma \\ & \quad \quad \quad 1 \leq i \leq l \\ & = \infty \end{aligned}$$

WRONG  
MARGIN

$$\begin{aligned} & \text{MAX}_{w,b,\gamma} \quad \gamma \leftarrow \text{GEOMETRIC} \\ & \text{SUBJECT TO} \quad y_i \left( \frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|} \right) \geq \gamma \quad (*) \\ & \quad \quad \quad 1 \leq i \leq l \end{aligned}$$

GOOD DEF. BUT HARD TO OPTIMIZE

# 6

## Hyperplanes with functional margin 1 known as canonical hyperplanes

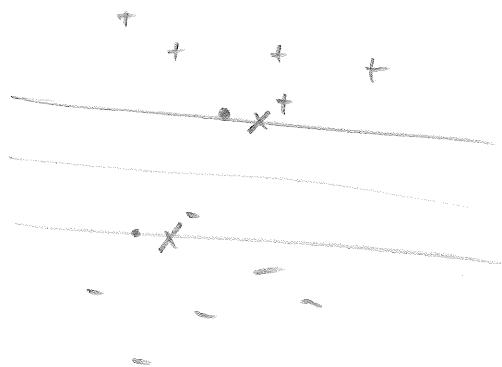
Find  $w, b$

$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1$$

$$1 \leq i \leq l$$

choose  $x^+, x^-$

$$\begin{aligned} \text{s.t. } w \cdot x^+ + b &= 1 \\ w \cdot x^- + b &= -1 \end{aligned}$$



Geometric margin of  $x^+$

$$\frac{w}{\|w\|^2} \cdot x^+ + \frac{b}{\|w\|^2} = \frac{1}{\|w\|^2} \underbrace{(x^+ + b)}_{+1}$$

Geometric margin of  $x^-$

$$\frac{w}{\|w\|^2} \cdot x^- + \frac{b}{\|w\|^2} = \frac{1}{\|w\|^2} \underbrace{(x^- + b)}_{-1}$$

## EQUIVALENT OPTIMIZATION PROBLEMS

$$\underset{w, b}{\text{MAX}} \quad \frac{1}{\|w\|^2}$$

$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1$$

MAX MARGIN  
HYPERPLANE

$$\underset{w, b}{\text{MIN}} \|w\|^2$$

$$\text{s.t. } y_i(w \cdot x_i + b) \geq 1$$

REALIZES MAX MARGIN  
HYPERPLANE W. MARGIN  $\frac{1}{\|w\|^2}$

## DUAL OPTIMIZATION PROBLEM?

$$L(w, b, \alpha) = \frac{1}{2} (w \cdot w) + \sum_{i=1}^l \alpha_i [y_i ((w \cdot \bar{x}_i + b) - 1)]$$

$$\frac{\partial L}{\partial w} = w + \sum_{i=1}^l \alpha_i y_i \bar{x}_i = 0$$

$$w = \sum_i y_i \alpha_i \bar{x}_i$$

LINEAR COMB. OF EXAMPLES

$$\frac{\partial L}{\partial b} = \sum_i y_i \alpha_i = 0$$

PLUG IN:

$$L = \frac{1}{2} \left( \sum_i y_i \alpha_i \bar{x}_i \right) \cdot \left( \sum_j y_j \alpha_j \bar{x}_j \right) -$$

$$+ \sum_i \alpha_i y_i \left( \sum_j \alpha_j y_j \bar{x}_j \cdot \bar{x}_i \right) - \underbrace{\sum_i \alpha_i y_i b}_{0} + \sum_i \alpha_i$$

$$= -\frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j \bar{x}_i \cdot \bar{x}_j + \sum_i \alpha_i$$

KEY: - ONLY DOT PRODUCTS MATTER

- KERNELIZABLE !!!

- KEY OBSERVATION OF ORIGINAL SUM PAPER

## DUAL PROBLEM

$$\text{MAXIMISE} \quad W(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y_i y_j \alpha_i \alpha_j x_i \cdot x_j$$

$$\text{SUBJECT TO} \quad \sum_i y_i \alpha_i = 0$$

$$\alpha_i \geq 0$$

IF  $\alpha^*$  IS SOLUTION TO DUAL

$$\text{THEN } w^* = \sum_i y_i \alpha_i^* x_i$$

$$b^* = \frac{\max_{y_i=-1} w^* \cdot x_i + \min_{y_i=+1} w^* \cdot x_i}{2}$$

## GEOMETRIC MARGIN

$$\frac{1}{\|w^*\|_2}$$

## KURESHI KUHN-TUCKER CONDITION

$$\alpha_i^* [ y_i (\underbrace{w^* \cdot x_i + b^*}_{\text{JFF MARGIN}}) - 1 ] = 0$$

$\nearrow$

0 JFF  
 MARGIN  
 > 1

> 0 JFF     "     > 1

## SUPPORT VECTORS :

ALL  $x_i$  FOR WHICH  $\alpha_i > 0$

## REPRESENTING MAX. MARGIN HYPERPLANE IN DUAL

$$\begin{aligned}
 w^* \cdot x + b^* &= \underbrace{\sum_{i=1}^l y_i \alpha_i^* x_i \cdot x}_{w^*} + b^* \\
 &= \sum_{i \in SV} y_i \alpha_i^* x_i \cdot x + b^*
 \end{aligned}$$

$$\text{FOR } j \in SV : \quad w^* \cdot x_j + b^* = y_j \quad (*)$$

$$\begin{aligned}
 \text{ALSO : } w^* \cdot w^* &= \sum_{i,j} y_i y_j \alpha_i^* \alpha_j^* x_i \cdot x_j \\
 &= \sum_{j \in SV} \alpha_j^* y_j - \underbrace{\sum_{i \in SV} y_i \alpha_i^* x_i \cdot x_j}_{w^*} \\
 (*) &= \sum_{j \in SV} \alpha_j^* y_j (y_j - b^*)
 \end{aligned}$$

$$y_j \in \{+1, -1\}$$

$$= \sum_{j \in SV} \alpha_j^* (1 - y_j b^*)$$

$$= \sum_{j \in SV} \frac{\alpha_j^*}{\gamma_0} - \underbrace{\sum_{j \in SV} \alpha_j^* y_j b^*}_0$$

$$\text{MARGIN : } \frac{1}{\|w^*\|_2} = \frac{1}{\sqrt{\sum_j \alpha_j^*}}$$

## KERNEL TRICK

$x_i$  EXPANDED TO  $\phi(x_i)$

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

$$w^* \cdot \phi(x) + b^* = \sum_{i \in SV} y_i x_i^* \phi(x_i) \cdot \phi(x) + b^*$$

↑  
ONE  
WEIGHT  
PER  
FEATURE

$$= \sum_{i \in SV} y_i x_i^* k(x_i, x) + b^*$$

↑

ONE WEIGHT  
PER SV

$|SV|$  USUALLY  $\ll$  DIMENSION OF FEATURE SPACE

GENERALIZATION ERROR GOOD IF

MARGIN LARGE

OR  $|SV|$  SMALL

## NON-SEPARABLE CASE

$$\underset{w, \xi}{\text{MINIMIZE}} \quad w \cdot w + C \sum_i \xi_i$$

$$\text{SUBJECT TO} \quad y_i (w \cdot x_i + b) \geq 1 - \xi_i$$

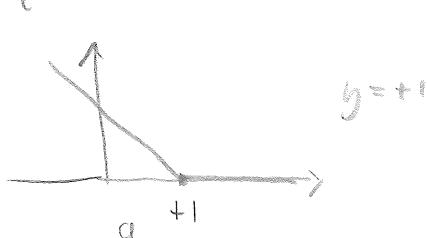
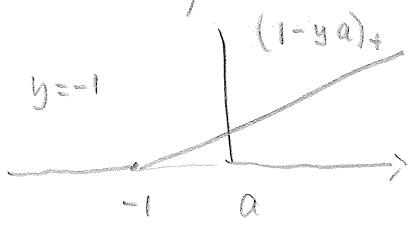
$$\xi_i \geq 0$$

ONE SLACK VARIABLE PER EXAMPLE

COST IS LINEAR IN SLACK

ALTERNATE DEF.

$$\underset{w, \xi}{\text{MINIMIZE}} \quad w \cdot w + C \sum_i (1 - y_i (\underbrace{w \cdot x_i + b}_{\alpha_i}))_+$$



(LINEAR) HINGE LOSS

DUAL

$$\text{MAXIMIZE: } \sum \alpha_i - \frac{1}{2} \sum_{ij} y_i y_j \alpha_i \alpha_j x_i \cdot x_j$$

$$\text{SUBJECT TO: } \sum y_i \alpha_i = 0$$

$$\alpha_i \in [0, C]$$

↑  
NEW

BEFORE  $C = \infty$

## QUADRATIC HINGE

$$\text{MINIMIZE: } w \cdot w + C \sum_i \xi_i^2$$

$$\text{SUBJ. TO} \quad y_i (w_i \cdot x_i + b) \geq 1 - \xi_i \\ \xi_i \geq 0$$

ALTERNATE FORMULATION:

$$\text{MINIMIZE } w \cdot w + C \sum_i ((1 - y_i (w_i \cdot x_i + b))^2)_+$$

## DUAL

$$\text{MAXIMIZE: } \sum_i \alpha_i - \frac{1}{2} \sum_{ij} y_i y_j \alpha_i \alpha_j x_i \cdot x_j \\ - \frac{1}{2C} \sum_i \alpha_i^2$$

$$\text{SUBJ TO: } \sum_i y_i \alpha_i = 0 \\ \alpha_i \geq 0$$

OTHER LOSSES  
PEOPLE INVESTI-  
GATED

15

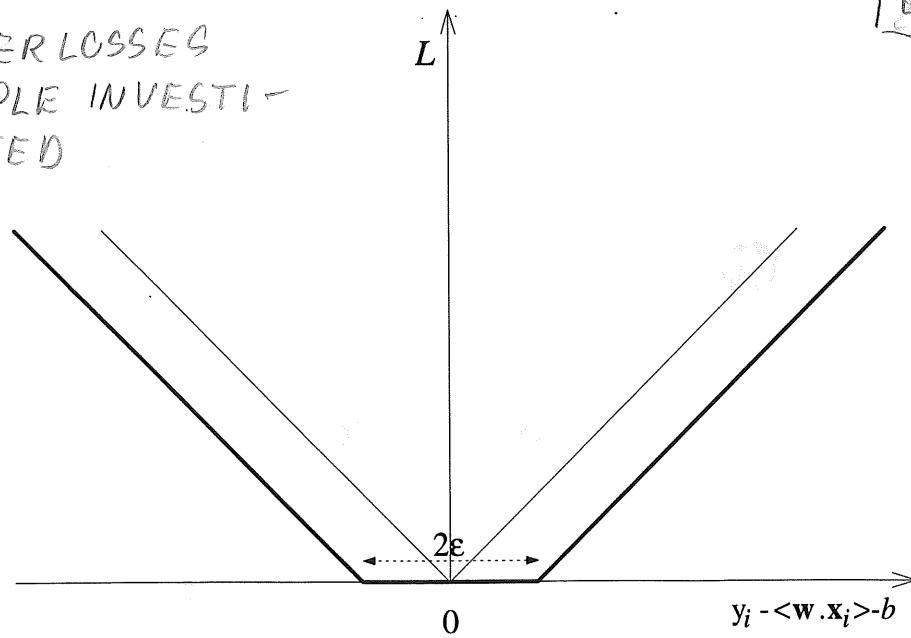


Figure 6.6: The linear  $\epsilon$ -insensitive loss for zero and non-zero  $\epsilon$

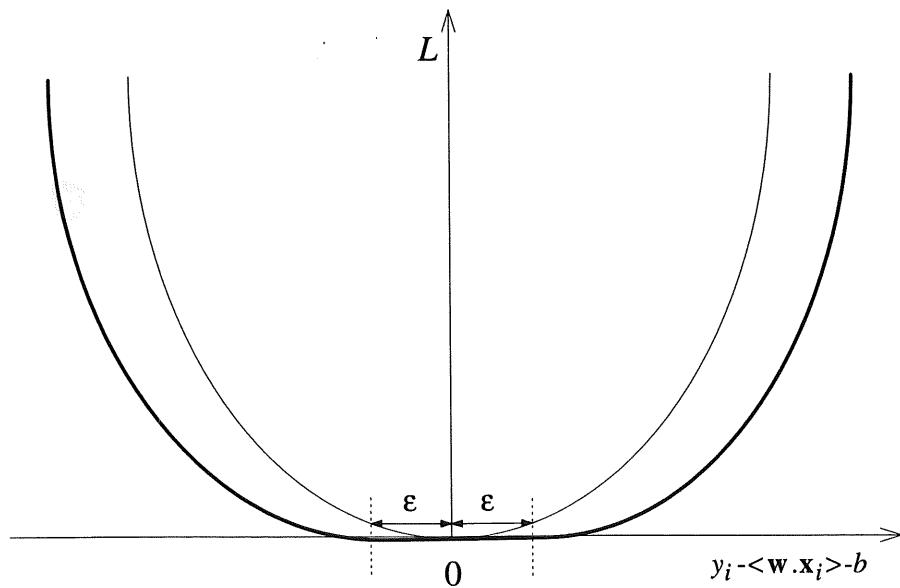


Figure 6.7: The quadratic  $\epsilon$ -insensitive loss for zero and non-zero  $\epsilon$

## KERNEL METHODS

WHEN APPLICABLE ?

HYPOTHESIS MUST BE DEFINED I.T.O.

DOT PRODUCT

$$\hat{y} = h_w(x) = f(w \cdot x)$$

PREDICTION ON NEW X

WHERE W IS LINEAR COMBINATION  
OF INSTANCES IN TRAINING SET

$$\text{i.e.: } w = \sum_{q=1}^t \alpha_q x_q$$

↑ CALLED LINEAR FORM

EXAMPLES :

WIDROW HOFF

PERCEPTRON

BACKPROP

SUPPORT VECTOR MACHINES

} GD FAMILY

W. IN LINEAR FORM WHEN

- PARAMETER DIVERGENCE IS  $\|w\|_2^2$

- LOSS DEPENDS ON DOT PRODUCT

$$w_{t+1} = \underset{w}{\operatorname{ARGINF}} \left( \|w - w_t\|^2 + \gamma \sum y_t (w \cdot x_t) \right)$$

$$\Rightarrow w_{t+1} = \sum_{q=1}^t \alpha_q x_q \quad \text{LINEAR FORM}$$

GD FAMILY

MAIN OTHER FAMILY

$$w_{t+1} = \underset{w \geq 0}{\operatorname{ARGINF}} \left( \sum w_i \ln \frac{w_i}{w_{t,i}} + w_{t,i} - w_i + \gamma \sum y_t (w \cdot x_t) \right)$$

$$\ln w_{t,i} = \sum_{q=1}^t \gamma_q x_{q,i}$$

$$w_{t,i} = e$$

EXPONENTIAL FORM

WINNOW,  
EGN FAMILY

ASSUME PREDICTION  $\hat{y}_t$  IS LINEAR ( $\hat{y}_t = w_t \cdot x_t$ )  
 OR LINEAR THRESHOLD FUNCTION ( $\hat{y}_t = (w_t \cdot x_t > \theta)$ )

LINEAR FORM GOOD BECAUSE  
 AFTER EXPANDING INSTANCES  
 STILL EFFICIENT

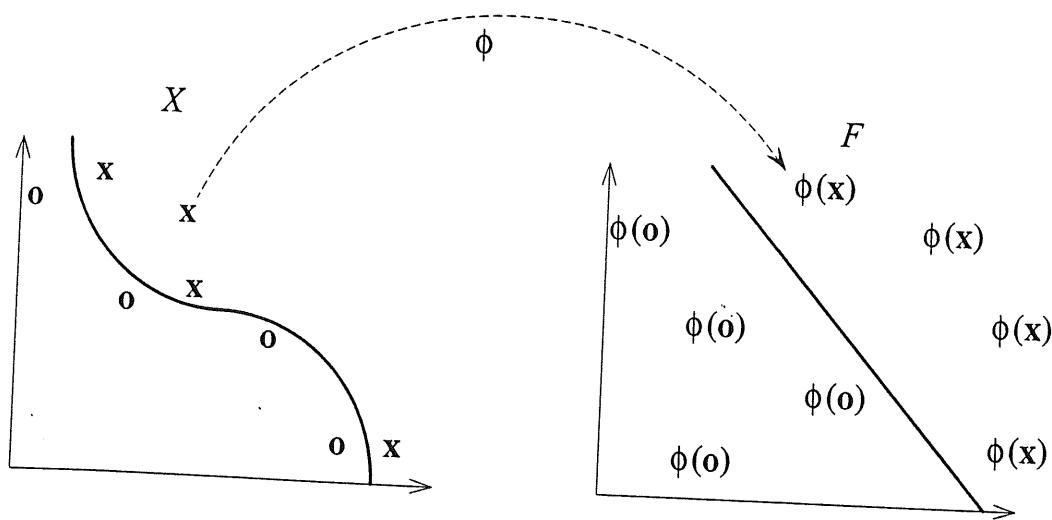


Figure 3.1: A feature map can simplify the classification task

$$\bar{x} = (x_1, \dots, x_n) \rightarrow \phi(\bar{x}) = (\phi_1(x), \dots, \phi_N(x))$$

$x \in X$

INPUT  
SPACE

$\{\phi(x) : x \in X\}$

FEATURE  
SPACE

FOR EXAMPLE :  $\phi(x_1, x_2) = (x_1^2, x_2^2, x_1 x_2)$

## REASONS FOR EXPANSION :

- LINEAR MODELS IN FEATURE SPACE ARE CONVENIENT NON-LINEAR MODELS OF INPUT SPACE

- DIMENSIONALITY REDUCTION  
(I.E.  $N < n$ )

- $\underbrace{\left( \sum_{q=1}^{t-1} \alpha_q^t \phi(x_q) \right)}_{\text{WEIGHT VECTOR}} \circ \phi(x_t)$

$w_t$  IN LINEAR FORM

$$= \sum_{q=1}^{t-1} \alpha_q^t \phi(x_q) \cdot \phi(x_t)$$

COMPUTE DOT PRODUCTS WITH PAST EXAMPLES

OFTEN EFFICIENT COMPUTATION OF  $\phi(x) \cdot \phi(z)$  EVEN WHEN

$$N \gg n$$

$\phi(x)$ 

$$(x_1, x_2, x_3) \rightarrow (1, x_1, x_2, x_3, x_1x_2, x_1x_3, x_2x_3, x_1x_2x_3)$$

8 FEATURES

Repeat!

$$\begin{aligned}
 \phi(x) \cdot \phi(z) &= 1 + x_1 z_1 + x_2 z_2 + x_3 z_3 + \\
 &\quad + x_1 x_2 z_1 z_2 + x_1 x_3 z_1 z_3 + x_2 x_3 z_2 z_3 \\
 &\quad + x_1 x_2 x_3 z_1 z_2 z_3 \\
 &= (1 + x_1 z_1) \cdot (1 + x_2 z_2) \cdot (1 + x_3 z_3)
 \end{aligned}$$

8 TERMS : 2 2 2 2

$$\begin{array}{ccc}
 & \phi(x) & \\
 \times & & \\
 (x_1, x_2, \dots, x_n) & (1, x_i, x_i x_j, x_i x_j x_k, \dots) & \\
 & \text{if } i \neq j \text{ all } \neq & \\
 & & 2^n \text{ MONOMIALS}
 \end{array}$$

$$\begin{aligned}
 \phi(x) \cdot \phi(z) &= \sum_{I \subseteq \{1, \dots, n\}} \prod_{i \in I} x_i \prod_{i \in I} z_i \\
 &= \prod_{i=1}^n (1 + x_i z_i)
 \end{aligned}$$

$O(n)$  TIME

EFFICIENCY:

$$\left( \sum_{q=1}^{t-1} \alpha_q^t \phi(q) \right) \cdot \phi(x_t)$$

$\underbrace{\phantom{\sum_{q=1}^{t-1} \alpha_q^t \phi(q)}}$

$w_t$

DIMENSION  $2^n$

SEEMINGLY TIME  $O(2^n)$

$$\begin{aligned}
 &= \sum_{q=1}^{t-1} \alpha_q^t \underbrace{\phi(q) \cdot \phi(x_t)}_{\text{TIME } O(n)} \\
 &\quad \underbrace{\phantom{\sum_{q=1}^{t-1} \alpha_q^t \phi(q) \cdot \phi(x_t)}}_{\text{TIME } O(t^n)}
 \end{aligned}$$

## MORE EXAMPLES

$$(x_1, \dots, x_n) \rightarrow \left( \begin{array}{c} x_i x_j \\ \scriptstyle 1 \leq i \leq n \\ \scriptstyle 1 \leq j \leq n \end{array} \right)$$

WITH REPEATS

$$\bar{x}$$

$$\phi(\bar{x})$$

$$\phi(\bar{x}) \cdot \phi(\bar{z}) = \sum_{\substack{(i,j) \\ (i,j) = (1,1)}}^{(n,n)} (x_i x_j)(z_i z_j)$$

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

$$= \left( \sum_{i=1}^n x_i z_i \right)^2$$

$$= (\bar{x} \cdot \bar{z})^2$$

## IMPLICIT MAPPING INTO FEATURE SPACE

A KERNEL IS A FUNCTION  $K$  S.T. FOR ALL  $x, z$

$$k(x, z) = \phi(x) \cdot \phi(z)$$

implicit

WHERE  $\phi$  IS A MAPPING FROM THE INPUT SPACE  $X$  TO A FEATURE SPACE  $F$  WHICH HAS AN INNER PRODUCT

## EXAMPLES SO FAR

$$k(x, z) = \prod_{i=1}^n (1 + x_i z_i)$$

$t$  training examples  
 $w = \sum_t \alpha_t x_t$

hypothesis repr. by  
 $t \alpha_t$ 's

$w \cdot x = t$  kernel  
 $\uparrow$  computations  
 $n \times n$

$$k(x, z) = (x \cdot z)^2$$

ALSO:

$$k(x, z) = (x \cdot z + c)^d$$

## PROPERTIES OF KERNEL FUNCS :

$$\left. \begin{aligned} k(x, z) &= \phi(x) \cdot \phi(z) \\ &= \phi(z) \cdot \phi(x) \\ &= k(z, x) \end{aligned} \right| \quad \begin{aligned} k(x, z)^2 &= (\phi(x) \cdot \phi(z))^2 \\ &\leq \|\phi(x)\|^2 \|\phi(z)\|^2 \\ &= (\phi(x) \cdot \phi(x)) (\phi(z) \cdot \phi(z)) \\ &= k(x, x) k(z, z) \end{aligned}$$

SYMMETRY

CAUCHY-SCHWARTZ

## CHARACTERIZATION OF KERNELS

MERCER'S TH (FINITE CASE) :

LET  $X = \{x_1, \dots, x_n\}$  BE A FINITE INPUT SPACE. THEN  $k(x_i, z)$  IS A KERNEL FUNCTION IFF

$$K = (k(x_i, x_j))_{i,j=1}^n$$

IS SYMMETRIC POSITIVE DEFINITE  
(HAS NON-NEG. EIGEN VALUES).

THERE ARE CONTINUOUS VERSIONS OF  
THIS THEOREM

WHAT ARE GOOD KERNELS  
FOR A GIVEN PROBLEM

## Special kernel functions

$$k(x, z) = h(x - z)$$

one translation invariant

Example:  $k(u) = \sum_{n=1}^{\infty} a_n \cos(nu)$

$$\begin{aligned} k(x-z) &= a_0 + \sum_{n=1}^{\infty} a_n \sin nx \sin nz \\ &\quad + \sum_{m=1}^{\infty} a_m \cos(nx) \cos(mz) \end{aligned}$$

$$\left( \phi_i(x) \right)_{i=0}^{\infty} = (1, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots, \sin(nx), \cos(nx))$$

FOURIER COEFFICIENTS

orthogonal features

$$\int_0^{2\pi} \sin(x) \cdot \cos(3x) = 0$$

$$k(x, z) = e^{-\frac{(x-z)^2}{2}}$$

rotation invar. as well

## MAKING KERNELS FROM KERNELS

**Proposition 3.12** Let  $K_1$  and  $K_2$  be kernels over  $X \times X$ ,  $X \subseteq \mathbb{R}^n$ ,  $a \in \mathbb{R}^+$ ,  $f(\cdot)$  a real-valued function on  $X$ ,

$$\phi : X \longrightarrow \mathbb{R}^m$$

with  $K_3$  a kernel over  $\mathbb{R}^m \times \mathbb{R}^m$ , and  $\mathbf{K}$  a symmetric positive semi-definite  $n \times n$  matrix. Then the following functions are kernels:

1.  $K(\mathbf{x}, \mathbf{z}) = K_1(\mathbf{x}, \mathbf{z}) + K_2(\mathbf{x}, \mathbf{z})$ ,
2.  $K(\mathbf{x}, \mathbf{z}) = aK_1(\mathbf{x}, \mathbf{z})$ ,
3.  $K(\mathbf{x}, \mathbf{z}) = K_1(\mathbf{x}, \mathbf{z})K_2(\mathbf{x}, \mathbf{z})$ ,
4.  $K(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})f(\mathbf{z})$ ,
5.  $K(\mathbf{x}, \mathbf{z}) = K_3(\phi(\mathbf{x}), \phi(\mathbf{z}))$ ,
6.  $K(\mathbf{x}, \mathbf{z}) = \mathbf{x}'\mathbf{B}\mathbf{z}$ .

**Corollary 3.13** Let  $K_1(\mathbf{x}, \mathbf{z})$  be a kernel over  $X \times X$ ,  $\mathbf{x}, \mathbf{z} \in X$ , and  $p(x)$  a polynomial with positive coefficients. Then the following functions are also kernels:

1.  $K(\mathbf{x}, \mathbf{z}) = p(K_1(\mathbf{x}, \mathbf{z}))$ ,
2.  $K(\mathbf{x}, \mathbf{z}) = \exp(K_1(\mathbf{x}, \mathbf{z}))$ ,
3.  $K(\mathbf{x}, \mathbf{z}) = \exp(-\|\mathbf{x} - \mathbf{z}\|^2 / \sigma^2)$ .