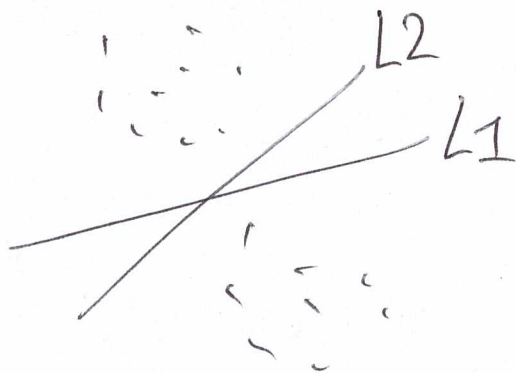


# Lecture 8: Support Vector Machines

①

## \* Linear classifier



$$D = \{ (x_i, y_i) \}$$

$\nearrow$   $x_i \in \mathbb{R}^n$        $\nwarrow$   $y_i \in \{ \pm 1 \}$

$$F(x) = \langle w, x \rangle - b$$

$\nearrow$  weight vector       $\nwarrow$  bias

$$\begin{cases}
 \langle w, x_i \rangle - b > 0 & \text{if } y_i = +1 \\
 \langle w, x_i \rangle - b < 0 & \text{if } y_i = -1
 \end{cases}$$

$$\rightarrow y_i (\langle w, x_i \rangle - b) > 0 \quad \forall (x_i, y_i) \in D$$

If possible:  $D$  is linearly separable

② Margin = distance hyperplane to closest data points

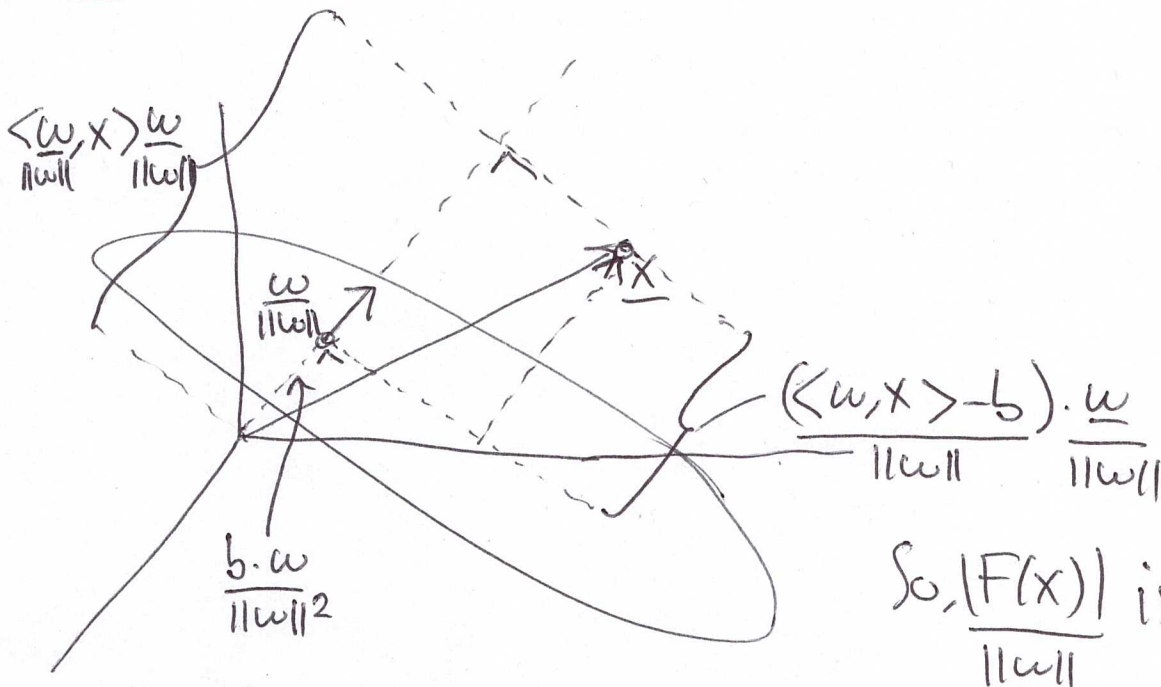
$$y_i (\langle w, x_i \rangle - b) \geq 1 \quad \forall (x_i, y_i) \in D$$

$\downarrow$   $= |F(x_i)|$        $\uparrow$  After rescaling       $\nearrow$  Some  $x_i$  will meet equality

Distance  $x_i$  to hyperplane:  $\frac{|F(x_i)|}{\|w\|}$

Margin =  $\frac{1}{\|w\|}$

$\updownarrow$   
 support vectors



So,  $\frac{|F(x)|}{\|w\|}$  is a proper distance measure.

Training problem:

Primal  $\left\{ \begin{array}{l} \text{minimize } Q(w) = \frac{1}{2} \|w\|^2 \\ \text{st. } y_i (\langle w, x_i \rangle - b) \geq 1 \quad \forall (x_i, y_i) \in D \end{array} \right.$  (for convenience)

- \*  $Q$  is convex in  $w$
- \* constraints linear in  $w$

Lagrange multipliers:

$$J(w, b, \alpha) = \frac{1}{2} \langle w, w \rangle - \sum_{i=1}^m \alpha_i \{ y_i (\langle w, x_i \rangle - b) - 1 \}$$

Minimize w.r.t.  $w$  and  $b \iff \frac{\partial J}{\partial w} = 0 \quad \& \quad \frac{\partial J}{\partial b} = 0$   
 Maximize w.r.t.  $\alpha \geq 0 \quad J^{\text{convex}} \iff \frac{\partial J}{\partial w} \updownarrow \quad \frac{\partial J}{\partial b} \updownarrow$

$w = \sum \alpha_i y_i x_i \quad \sum \alpha_i y_i = 0$   
 (a) (b)

# Dual Problem

- \* Same optimal value as primal problem
- \*  $\alpha_i$ 's provide optimal solution

$$J(w, b, \alpha) = \underbrace{\frac{1}{2} \langle w, w \rangle - \sum \alpha_i y_i \langle w, x_i \rangle}_{by (a)} - \underbrace{b \sum \alpha_i y_i + \sum \alpha_i}_{= 0 by (b)}$$

$$\langle w, w \rangle \stackrel{by (a)}{=} \sum \alpha_i y_i \langle w, x_i \rangle \stackrel{by (a)}{=} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

$$J(w, b, \alpha) = Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

Dual } maximize  $Q(\alpha)$   
 st.  $\sum \alpha_i y_i = 0, \alpha \geq 0$

\* Just training data  $\mathbb{P}$

\*  $Q(\alpha)$  only depends on  $\langle x_i, x_j \rangle$   $\mathbb{P}$

Optimum  $\alpha_i^* \Rightarrow w^* = \sum_i \alpha_i^* y_i x_i$

Kuhn-Tucker  $\rightsquigarrow x_i$  is a support vector iff  $\alpha_i^* \neq 0$   
 $\downarrow$   
 $b^* = 1 - \langle w^*, x_i \rangle$

Not linearly separable

Allow mislabeled data → use slack variables  $\zeta_i$

measure degree of misclassification

Primal

Minimize  $Q(w, b, \zeta) = \frac{1}{2} \|w\|^2 + C \cdot \sum \zeta_i$

Trade-off:

- margin size
- amount of error in training

s.t.  $y_i (\langle w, x_i \rangle - b) \geq 1 - \zeta_i \quad \forall (x_i, y_i) \in D$   
 $\zeta_i \geq 0$

Dual

Maximize  $Q(\alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$   
(mistake in book)

s.t.  $\sum_i \alpha_i y_i = 0, \quad C \geq \alpha > 0$

(In the proof: use Lagrange variables  $\mu_i$  to enforce non-negativity of slack variables  $\zeta_i$ )

$\zeta_i = 0$  if  $\alpha_i < C$   
 $\zeta_i \geq 0$  if  $\alpha_i = C$  )

# Slides [6-11] + [15-21] tutorial

## Kernel trick

Non-linear SVMs: 1)  $x_i \rightarrow \varphi(x_i)$  in higher dim. space

2)  $\{(\varphi(x_i), y_i)\}$  can be linearly separated

$$Q(x) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underbrace{\langle \varphi(x_i), \varphi(x_j) \rangle}_{K(x_i, x_j)}$$

$$\exists \varphi \quad K(u, v) = \langle \varphi(u), \varphi(v) \rangle$$

$$\forall \varphi \quad \left[ \int \varphi(x)^2 dx \leq 0 \Rightarrow \int K(u, v) \varphi(u) \varphi(v) du dv \leq 0 \right]$$

We do not need to know  $\varphi$

We do not need to compute a large inner product



$$F(x) = \sum_i \alpha_i^* y_i k(x_i, x) - b$$

## Slides [26-28] tutorial

SVM Regression: estimate a function that maps input objects  $x_i \in \mathbb{R}^n$

→ real number  $y_i \in \mathbb{R}$ , not necessarily in  $\{-1, 1\}$   
based on training data

$$F(x) = \langle w, x \rangle - b$$

How to best-fit the data?

$$|y_i - F(x_i)| \leq \epsilon, \text{ margin } \frac{1}{\|w\|}$$

$$\begin{cases} \text{Minimize } L(w) = \frac{1}{2} \|w\|^2 \\ \text{s.t. } |y_i - F(x_i)| \leq \epsilon \quad \forall (x_i, y_i) \in D \end{cases}$$

Allow noise in training data

$$\text{Minimize } L(w, \zeta, \hat{\zeta}) = \frac{1}{2} \|w\|^2 + C \cdot \left( \sum_i \zeta_i^2 + \sum_i \hat{\zeta}_i^2 \right)$$

$$\text{s.t. } \begin{aligned} y_i - (\langle w, x_i \rangle - b) &\leq \epsilon + \zeta_i \\ -y_i + (\langle w, x_i \rangle - b) &\leq \epsilon + \hat{\zeta}_i \end{aligned}$$

$$\zeta_i, \hat{\zeta}_i \geq 0$$

etc.

SVM Ranking : learning a ranking/preference function

$\{(x_i, y_i)\}$  :  $x_i$  is preferred to  $x_j$  ( $x_i \succ x_j$ )  
if  
 $y_i < y_j$

Output score :  $F(x_i) > F(x_j)$  for any  $x_i \succ x_j$   
(not class)

Linear ranking function  $\langle w, x_i \rangle > \langle w, x_j \rangle$

Approximate solution using slack variables:

$$\text{minimize } \frac{1}{2} \langle w, w \rangle + C \sum \zeta_{ij}$$

$$\text{s.t. } \forall_{y_i < y_j} \quad \langle w, x_i \rangle \geq \langle w, x_j \rangle + 1 - \zeta_{ij}$$
$$\zeta_{ij} \geq 0$$

etc  
(dual)

$$F(z) = \sum_{i,j} \alpha_{ij}^* K(x_i - x_j, z)$$

# 1-Norm ranking SVM

(P)

$$F(x_u) > F(x_v) \Rightarrow \sum_{(i,j) \in P} \alpha_{ij} \langle (x_i - x_j), x_u \rangle > \sum_{(i,j) \in P} \alpha_{ij} \langle (x_i - x_j), x_v \rangle$$

$(u,v) \in P$

$$P = \{(i,j) : y_i < y_j\}$$

$$L(\alpha, \beta) = \sum_{(i,j) \in P} \alpha_{ij} + C \sum_{(i,j) \in P} \beta_{ij}$$

$$\forall y_u < y_v \quad \sum_{(i,j) \in P} \alpha_{ij} \langle (x_i - x_j), (x_u - x_v) \rangle \geq 1 - \beta_{uv}$$

$$\alpha_{ij} \geq 0, \beta_{u,v} \geq 0$$