

# SUPPORT VECTOR MACHINES

①

- Biblios :
- Bishop, "Pattern Recognition and Machine Learning", Cap 7.
  - Schölkopf & Smola, "Learning with Kernels", cap 7.

• limitante de los métodos de kernel: (i) si se requiere evaluar  $k(\underline{x}_n, \underline{x}_m)$  para todo par de muestras de entrenamiento, puede volverse inrealizable; (ii) testing o clasificación de  $\underline{x}$ : costo computacional de evaluar  $k(\underline{x}, \underline{x}_n)$  para todos  $\underline{x}_n$  puede ser excesivo.

• SVM: Veremos que

- la clasificación de nuevas muestras depende sólo de la evaluación de funciones kernel en un subconjunto extremadamente reducido del conjunto de entrenamiento.

- los parámetros del clasificador se determinan fácilmente, resolviendo un problema de optimización convexa.

## 1. PROBLEMA DE DOS CLASES

•  $\{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N\}$  muestras de entrenamiento,  $\underline{x}_n \in \mathcal{X}$   
 $\underline{x}_n \leftrightarrow$  etiqueta  $t_n \in \{-1, 1\}$ ,  $n=1, \dots, N$

•  $\phi(\cdot)$  mapeo no lineal a un espacio de dimensión mayor:  $\phi: \mathcal{X} \rightarrow \mathcal{Y}$ ,  $\dim \mathcal{Y} > \dim \mathcal{X}$

Ec. de superficie de decisión  $\mathcal{Y}$ :

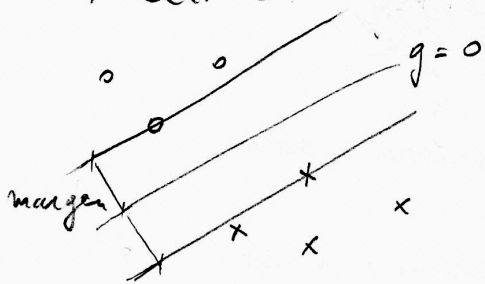
$$g(\underline{x}) = \underline{w}^T \phi(\underline{x}) + b, \quad \mathcal{Y} = \{ \underline{x} : g(\underline{x}) = 0 \}$$

En el espacio transformado  $\mathcal{Y}$ , la superficie de decisión es un hiperplano de vector normal  $\frac{\underline{w}}{\|\underline{w}\|}$  y sesgo  $b/\|\underline{w}\|$

• Una nueva muestra  $\underline{x}$  se clasifica según  $\text{signo}(g(\underline{x}))$

## 1.1. CLASES LINEALMENTE SEPARABLES (EN APRENDIZAJE)

- En este caso, existe al menos un par  $(\underline{w}, b)$  para el cual todas las muestras de entrenamiento están bien clasificadas, i.e.  $t_n g(\underline{x}_n) > 0 \quad \forall i=1, \dots, N$
- Entre los  $(\underline{w}, b)$  posibles, SVM elige aquel que corresponde al hiperplano separador con margen máximo.



Es intuitivo (y se puede demostrar [Vapnik]) que maximizar el margen minimiza el error de generalización del clasificador.

- Distancia de un punto  $\underline{x}$  al plano  $\{g(\underline{x})=0\}$ :  $\frac{|g(\underline{x})|}{\|\underline{w}\|}$
- Para cualquier hiperplano separador, puesto que estamos en el caso separable:
- $$|g(\underline{x}_n)| = t_n g(\underline{x}_n), \quad n=1, 2, \dots, N$$

Solución de máximo margen:

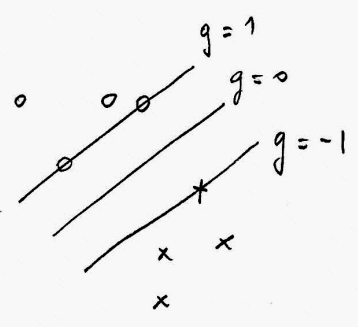
$$\arg \max_{\underline{w}, b} \left\{ \frac{1}{\|\underline{w}\|} \min_{n \in \{1, \dots, N\}} [t_n (\underline{w}^T \phi(\underline{x}_n) + b)] \right\}$$

Esta optimización es complicada, pero podemos encontrar una formulación equivalente sencilla de optimizar:

Obs: el escalado de  $\underline{w}$  y  $b$  por una misma constante no altera el cociente  $\frac{t_n g(\underline{x}_n)}{\|\underline{w}\|}$

$\Rightarrow$  Podemos definir la normalización tal que  $t_k (\underline{w}^T \phi(\underline{x}_k) + b) = 1$  para el (o los) puntos  $\underline{x}_k$  que realizan el margen (el más cercano a  $\{g(\underline{x})=0\}$ )

luego, tenemos que  $\forall n=1, 2, \dots, N, t_n (\underline{w}^T \phi(\underline{x}_n) + b) \geq 1$ .



- se llaman :
  - puntos activos : aquellos para los cuales  $t_n(w^T \phi(x_n) + b) = 1$
  - puntos inactivos : el resto ( $> 1$ )

• Siempre habrá al menos un punto activo, y al maximizar el margen habrá al menos dos

Tenemos entonces la siguiente formulación equivalente del problema de maximización del margen con clasificación (de entrenamiento) perfecta :

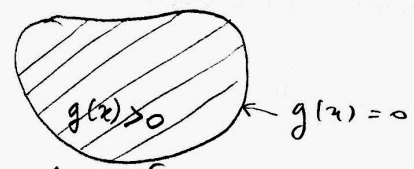
$$\begin{cases} \min_{w, b} \frac{1}{2} \|w\|^2 \\ \text{sujeto a } t_n (w^T \phi(x_n) + b) \geq 1, \quad n=1, \dots, N \end{cases}$$

Es un problema de programación cuadrática, fácil de resolver. Hay "solvers" que vienen en el toolbox de optimización de Matlab (ver también el toolbox no propietario CVX).

Pero el problema anterior (llamado primal) se puede llevar a una formulación aún más sencilla mediante las condiciones de Karush-Kuhn-Tucker (KKT) (llamado problema dual).

BREVE PARENTESIS : CONDICIONES KKT

$$\begin{cases} \min f(x) \\ \text{s.t. } g(x) \geq 0 \end{cases}$$

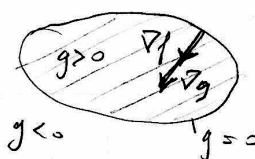


2 posibilidades para la solución :

- (i)  $x \in \{g(x) > 0\}$  : en este caso  $x$  es mínimo local de  $f \Rightarrow \nabla f(x) = 0$
- (ii)  $x \in \{g(x) = 0\}$  : en este caso debe ser  $\nabla f(x) \perp \{g=0\}$ , ya que de lo contrario existiría un vecino  $x'$  de  $x$ ,  $x' \in \{g=0\}$ , con  $f(x') < f(x)$ .



Como  $\nabla g \perp \{g=0\}$ ,  $\exists \lambda \neq 0$  t.f.  $\nabla f(x) = \lambda \nabla g(x)$ . luego, debe ser  $\lambda > 0$ , ya que si  $\lambda < 0$ , existiría un vecino  $x'$  de  $x$ ,  $x' \in \{g > 0\}$ , para el cual  $f(x') < f(x)$ . (4)



Definiendo el Lagrangeano

$Z(x, \lambda) = f(x) - \lambda g(x)$ , las 2 situaciones se resumen en una única, mediante las condiciones:

$$\frac{\partial Z}{\partial x} = 0, \quad \lambda \geq 0, \quad \lambda g(x) = 0 \quad (\text{a esta última se le llama condición kKT})$$

Obs 1: La solución  $(x^*, \lambda^*)$  es punto silla de  $Z(x, \lambda)$ :

$$Z(x^*, \lambda) \leq Z(x^*, \lambda^*) \leq Z(x, \lambda^*), \quad \forall (x, \lambda)$$

En efecto,  $Z(x^*, \lambda) \leq Z(x^*, \lambda^*)$

$$\Rightarrow \underbrace{f(x^*) - \lambda g(x^*)}_{\leq 0} \leq \underbrace{f(x^*) - \lambda^* g(x^*)}_{=0} \quad \checkmark$$

$Z(x^*, \lambda^*) \leq Z(x, \lambda^*)$

$$\Rightarrow \underbrace{f(x^*) - \lambda^* g(x^*)}_{=0} \leq f(x) - \lambda^* g(x)$$

$$\Rightarrow \underbrace{\lambda^* g(x)}_{\geq 0} \leq f(x) - f(x^*) \quad \checkmark$$

Obs 2: De obs 1; la optimización correspondiente es

$$\min_x \max_{\lambda} Z(x, \lambda).$$

Volviendo a nuestro problema:

$$Z(\underline{w}, b, \underline{a}) = \frac{1}{2} \|\underline{w}\|^2 - \sum_{n=1}^N a_n \left( \underbrace{t_n (\underline{w}^T \phi(\underline{x}_n) + b)}_{\phi^T(\underline{x}_n) \underline{w}} - 1 \right)$$

$$(P) \begin{cases} \min_{\underline{w}, b} \max_{\underline{a}} Z(\underline{w}, b, \underline{a}) \\ \text{con } a_n \geq 0, \quad n=1, 2, \dots, N \end{cases}$$

son las restricciones de positividad

$$\frac{\partial Z}{\partial \underline{w}} = 0 \Rightarrow \underline{w} = \sum_{n=1}^N a_n t_n \phi(\underline{x}_n) \quad (*)$$

$$\frac{\partial Z}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n t_n = 0 \quad (**)$$

Reemplazamos (\*) en  $Z(\underline{w}, b, \underline{a})$ :



(5)

$$\begin{aligned} \mathcal{L}(\underline{w}, b, \underline{a}) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \phi^T(\underline{x}_n) \phi(\underline{x}_m) \\ &\quad - \sum_{n=1}^N a_n t_n \phi^T(\underline{x}_n) \sum_{m=1}^N a_m t_m \phi(\underline{x}_m) \\ &\quad - b \underbrace{\sum_{n=1}^N a_n t_n}_{=0 \text{ (ec **)}} + \sum_{n=1}^N a_n \end{aligned}$$

$$\mathcal{L}(\underline{w}, b, \underline{a}) = \tilde{\mathcal{L}}(\underline{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m \underbrace{\phi^T(\underline{x}_n) \phi(\underline{x}_m)}_{k(\underline{x}_n, \underline{x}_m)}$$

El nuevo problema, llamado problema dual, queda:

$$(P') \begin{cases} \max_{\underline{a}} \tilde{\mathcal{L}}(\underline{a}) \\ \text{con } \underline{a} \geq 0 \end{cases}$$

De nuevo, este es un problema de programación cuadrática, pero todavía más sencillo de resolver ya que involucra únicamente a los multiplicadores de Lagrange.

• Clasificación de una nueva muestra  $\underline{x}$ :

se evalúa el signo de  $g(\underline{x}) = \underline{w}^T \phi(\underline{x}) + b$

$$= \sum_{n=1}^N a_n t_n k(\underline{x}_n, \underline{x}) + b$$

Veremos como calcular  $b$  más adelante.

• Vectores de soporte

De las condiciones KKT;

$$\begin{aligned} a_n &\geq 0 \\ t_n g(\underline{x}_n) - 1 &\geq 0 \\ a_n (t_n g(\underline{x}_n) - 1) &= 0 \end{aligned}$$

Tenemos dos tipos de muestras de entrenamiento, que verifican de una forma u otra la última condición:

(i)  $a_n = 0$

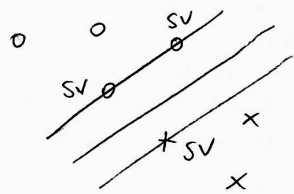
(ii)  $t_n g(\underline{x}_n) = 1$

los primeros no influyen en la clasificación de una nueva muestra:

(6)

$$g(\underline{x}_n) = \sum_{\substack{a_m \neq 0 \\ m=1, \dots, N}} a_m t_m k(\underline{x}_m, \underline{x}) + b$$

Los restantes se llaman vectores de soporte (SV), como verifican  $t_m g(\underline{x}_m) = 1$ , son muestras que caen sobre los hiperplanos de máximo margen.



Conclusión: en la etapa de testing, una vez que se ha entrenado el modelo, el proceso es rápido porque sólo se consideran los vectores de soporte ( $\#SV \ll N$  en general)

Cuanto vale el margen?

$$SV : t_m g(\underline{x}_m) = 1 \Rightarrow t_m \left( \sum_{\underline{x}_m \text{ SV}} a_m t_m k(\underline{x}_m, \underline{x}_m) + b \right) = 1$$

$$\Rightarrow \sum_{\underline{x}_m \text{ SV}} t_m \underbrace{\left( \sum_{\underline{x}_m \text{ SV}} a_m t_m k(\underline{x}_m, \underline{x}_m) + b \right)}_{=1} = \sum_{\underline{x}_m \text{ SV}} t_m$$

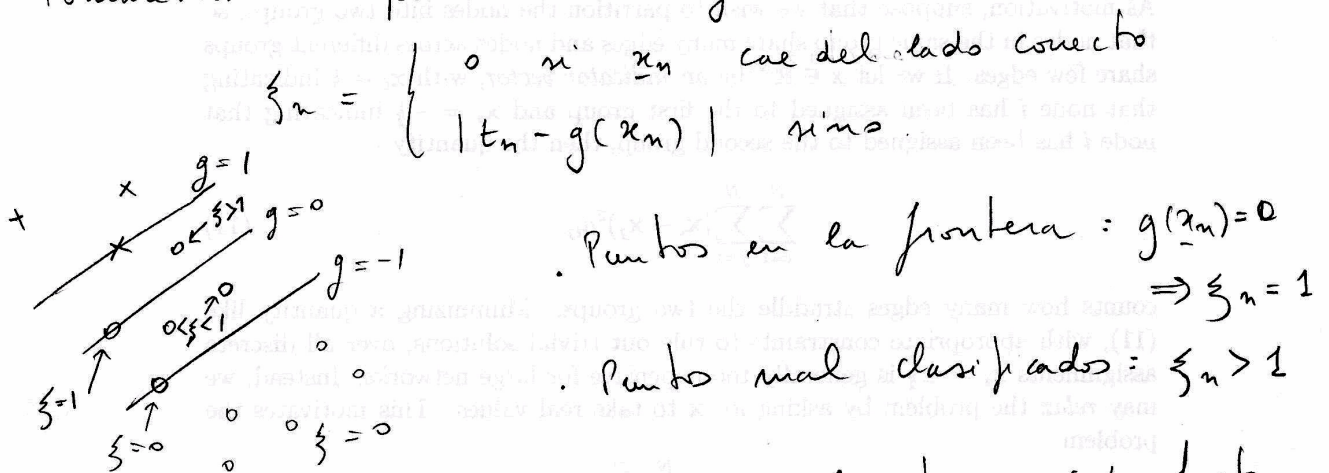
$$\Rightarrow b = \frac{1}{\#SV} \sum_{\underline{x}_m \text{ SV}} \left( t_m - \sum_{\underline{x}_m \text{ SV}} a_m t_m k(\underline{x}_m, \underline{x}_m) \right)$$

$\Rightarrow$   
 $\times t_m$   
 y sumo en los SV

## 1.2. CASO NO SEPARABLE LINEALMENTE (EN APRENDIZAJE) (7)

la idea es modificar la formulación anterior para autorizar que algunos puntos caigan del lado equivocado del hiperplano, pero penalizando esos errores, por ejemplo con una penalización que crezca con la distancia al hiperplano de decisión.

Tomaremos la penalización siguiente:



Reemplazamos la condición de clasificación perfecta

$$t_n (\underline{w}^T \phi(x_n) + b) \geq 1 \quad (\text{"Hard Margin"})$$

por la condición relajada o "Soft Margin":

$$\begin{cases} t_n (\underline{w}^T \phi(x_n) + b) \geq 1 - \xi_n, & n=1, \dots, N \\ \text{con } \xi_n \geq 0 \end{cases}$$

OBJETIVO: maximizar el margen, penalizando errores:

$$\begin{cases} \text{min}_{\underline{w}, b} \left\{ c \sum_{n=1}^N \xi_n + \frac{1}{2} \|\underline{w}\|^2 \right\}, & c > 0. \\ \text{sujeito a } \begin{cases} t_n (\underline{w}^T \phi(x_n) + b) \geq 1 - \xi_n, & n=1, \dots, N. \\ \xi_n \geq 0 \end{cases} \end{cases}$$

$c$ : constante de trade-off errores vs. margen.

$c \rightarrow 0$ : no autoriza errores, Hard Margin.

$c \downarrow$ : más errores, margen mayor.



El problema se resuelve de forma análoga:

$$\begin{aligned}
 \mathcal{L}(\underline{w}, b, \underline{a}, \underline{\mu}) &= \frac{1}{2} \|\underline{w}\|^2 + c \sum_{n=1}^N \xi_n - \sum a_n (\tan g(\underline{x}_n) - 1 + \xi_n) \\
 &\quad - \sum_{n=1}^N \mu_n \xi_n
 \end{aligned}$$

$a_n \geq 0, \mu_n \geq 0$  mult. de Lagrange

KKT :

$$\left\{ \begin{aligned}
 a_n &\geq 0 \\
 \tan g(\underline{x}_n) - 1 + \xi_n &\geq 0 \\
 a_n (\tan g(\underline{x}_n) - 1 + \xi_n) &= 0 \\
 \mu_n &\geq 0 \\
 \xi_n &\geq 0 \\
 \mu_n \xi_n &\geq 0
 \end{aligned} \right.$$

$$\frac{\partial \mathcal{L}}{\partial \underline{w}} = 0 \Rightarrow \underline{w} = \sum_{n=1}^N a_n \tan \phi(\underline{x}_n)$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N a_n \tan = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \Rightarrow a_n = c - \mu_n \Rightarrow \mu_n = c - a_n$$

Reemplazando  $\underline{w}, \mu_n$  en  $\mathcal{L}(\underline{w}, b, \underline{a}, \underline{\mu})$  con las ec. anteriores, los  $\xi_n$  y  $b$  se cancelan y obtenemos el Lagrangiano dual

$$\tilde{\mathcal{L}}(\underline{a}) = \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m \tan t_n t_m k(\underline{x}_n, \underline{x}_m)$$

$$\left. \begin{aligned}
 \text{con } a_n &\geq 0 \\
 \mu_n = c - a_n &\geq 0
 \end{aligned} \right\} \Rightarrow 0 \leq a_n \leq c$$

• Clasificación de nueva muestra  $\underline{x}$

(9)

según signo de  $g(\underline{x}) = \sum_{n=1}^N a_n t_n k(\underline{x}_n, \underline{x}) + b$

• Support vectors?

$$a_n > 0$$

$$t_n g(\underline{x}_n) = 1 - \xi_n$$

• Si  $a_n < C \Rightarrow \mu_n > 0 \Rightarrow \xi_n = 0$  : son SV dentro del lado correcto, fuera del margen.  
 $\uparrow$   
 $\mu_n \xi_n = 0$

• Si  $a_n = C \Rightarrow$  puede ser  $\mu_n = 0 \Rightarrow \begin{cases} \xi_n \leq 1 & \text{bien clasif.} \\ \xi_n > 1 & \text{mal clasificado} \end{cases}$

• Margen  $b$ ?

Support vectors:  $0 < a_n < C \Rightarrow \xi_n = 0 \Rightarrow t_n g(\underline{x}_n) = 1$

Exactamente de la misma forma que para el SVM Hard Margin, obtenemos:

$$b = \frac{1}{\#S} \sum_{\underline{x}_n \in S} \left( t_n - \sum_{\underline{x}_m \in S} a_m t_m k(\underline{x}_m, \underline{x}_n) \right)$$

donde  $S = \{ \underline{x}_n : 0 < a_n < C, n=1,2,\dots,N \}$   
Conjunto de SVs.

- El método Soft-Margin presentado aquí se conoce como C-SVM - Existen otras versiones soft margin (e.g.  $\nu$ -SVM, ver Schölkopf & Smola).

- Usualmente  $C$  y el  $\sigma$  los parámetros del kernel se estiman conjuntamente usando Grid Search y validación cruzada sobre el conjunto de entrenamiento.

## 2. EXTENSIÓN A M CLASES

(10)

### 2.1. Uno contra el resto

- Para cada una de las  $M$  clases, se construye un clasificador lineal usando como conjunto de entrenamiento la clase en cuestión por un lado, y todo el resto por el otro lado.

$$g^j(\underline{x}) = \sum_{n=1}^N a_n^j t_n^j k(\underline{x}_n, \underline{x}) + b^j, \quad j=1, \dots, M.$$

- Una muestra no etiquetada  $\underline{x}$ , se asigna a la clase  $\operatorname{argmax}_{j=1, \dots, M} g^j(\underline{x})$
- Desventaja: poco simétrico.

### 2.2. Clasificación por pares

- Para cada par de clases, se consideran únicamente los datos de esas dos clases y se entrena un clasificador lineal para separarlas  $\rightarrow M(M-1)/2$  clasificadores lineales
- Costo computacional: hay que resolver más problemas que para "uno contra todos", pero con menos muestras.
- Clasificación de una nueva muestra  $\underline{x}$ :
  - se evalúan los  $M(M-1)/2$  clasificadores, y  $\underline{x}$  se asigna a la clase más votada.

### 2.3. Error correcting output coding

- Se generan  $L$  particiones binarias de clases. Por ejemplo, para la base de dígitos  $0, 1, 2, \dots, 9$  USPS:

$$P_1 = (\{0, 1, \dots, 4\}, \{5, \dots, 9\}), \quad P_2 = (\{0, 1, 4, 5, 8\}, \{2, 3, 6, 7, 9\})$$
$$P_3 = (\{0, 2, 4, 6, 8\}, \{1, 3, 5, 7, 9\}) \dots$$

D: Hamming u otra distancia que además contemple cuán grande es cada error.

La idea es que a cada clase le corresponda un único código  $\underline{d} = (d_1, d_2, \dots, d_L)$ , con  $d_i \in \{-1, 1\}$  según en qué sub-clase de la partición  $P_i$  esté la clase. Luego se entrena un clasificador lineal para cada partición,  $f_1, \dots, f_L$ .

Clasificación:  $\underline{x} \rightarrow (f_1(\underline{x}), f_2(\underline{x}), \dots, f_L(\underline{x})) = F(\underline{x})$   
 $\underline{x}$  se asigna a  $\operatorname{argmin}_{j=1, \dots, M} D(F(\underline{x}), \underline{d}_j)$