

TEMA 7: CLASIFICADORES LINEALES ÓPTIMOS

MÁRGENES

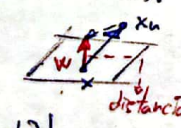
- + En casos lin. separables \rightarrow ¿Podemos preferir cierta h^* ?
 - * Si elegimos h : $\uparrow \uparrow$ margen \rightarrow intuitivamente, seremos menos susceptibles al "ruido" en los datos. $\rightarrow \downarrow \downarrow$ Error
- + Por tanto, buscaremos h : tenga el mayor margen posible.

CÁLCULO DEL HIPERPLANO SEPARADOR CON MÁRGEN MÁX.

PASOS PRELIMINARES

- + Normalización de w : Si multiplicamos todos los valores en w por una cte. " k " \rightarrow Nos quedamos con el mismo plano.
 - * Por simplicidad, buscaremos el w tal que $|w^T x_n| = 1$, siendo x_n el ejemplo más cercano al plano: $\forall x_n \in D \rightarrow |w^T x_n| \geq 1$.
- + Suponemos que los ejemplos son lin. separables: $y(w^T x_n) \geq 1$.
 $\text{signo}(y) = \text{signo}(w^T x_n)$
- + Cambiamos w_0 por $b \rightarrow$ Plano: $w^T x + b = 0$.

DISTANCIA PUNTO - PLANO. CÁLCULO DEL MÁRGEN.

- + El vector w es perpendicular al plano w .
 - * Si tomamos 2 puntos del plano: x' y x'' :
 $w^T x' + b = 0 \wedge w^T x'' + b = 0 \rightarrow w^T (x' - x'') = 0 \rightarrow w^T \perp \overrightarrow{x'x''}$
 $\overrightarrow{x'x''}$ vector sobre el plano
- + Distancia de x_n al plano:
 $\text{distancia}(w, x_n) = \left| \hat{w}^T (x_n - x) \right|$, $\hat{w}^T = \frac{w}{\|w\|}$ w normalizado ($\|w\|=1$).
 \rightarrow Proyección de $\overrightarrow{wx_n}$ en \hat{w} : 
- $\text{dist}(w, x_n) = \frac{1}{\|w\|} |w^T x_n - w^T x| = \frac{1}{\|w\|} |w^T x_n + b - (w^T x + b)|$
 $\text{min}(\text{dist}(w, x_n)) = \frac{1}{\|w\|} \rightarrow$ Margen de w .

OPTIMIZACIÓN DEL MÁRGEN

- + Tenemos que maximizar el margen $= \frac{1}{\|w\|} = \frac{1}{\sqrt{w^T w}}$.
- + Este problema se traduce en minimizar $w^T w$.
 - * Tenemos restricciones: $y_n (w^T x_n) \geq 1$ } - Ejemplos bien clasificados
 $(+b)$ } - Mínima valoración = 1.
 $\text{min}(y_n (w^T x_n + b)) = 1$

\rightarrow Podemos resolver el problema con "programación cuadrática".

- * P.C.: $\left\{ \begin{array}{l} - \text{Resuelve } \min_{u \in \mathbb{R}^n} \frac{1}{2} u^T Q u + p^T u \\ - \text{Con restricciones: } Au \geq c. \end{array} \right.$
- * En nuestro caso $\left\{ \begin{array}{l} - \text{Minimizamos } \frac{1}{2} w^T Q w + p^T w \\ \text{Hacemos que sea } \approx I \end{array} \right.$ No necesitamos esto, así que $p=0$.
 $- \text{Sujeto a: } A \cdot w \geq c \rightarrow c=1$
 $A = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \rightarrow Aw = y_n (w^T x_n) \rightarrow$ Incluye b .

\rightarrow Otra aproximación: Problema Dual de Lagrange.

- * $\mathcal{L}(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{n=1}^N \alpha_n (y_n (w^T x_n + b) - 1)$
 $\xrightarrow{\text{Lagrangiano}} \xrightarrow{\text{Los coeficientes de Lagrange.}}$
- * Minimizar $\mathcal{L}(w, b, \alpha)$ respecto a w y b , y después maximizar respecto a α resuelve nuestro problema.
 $(\alpha_n \geq 0)$
- 1. Minimizamos $\left\{ \begin{array}{l} \nabla_w \mathcal{L} = w - \sum_{n=1}^N \alpha_n y_n x_n = 0 \\ \nabla_b \mathcal{L} = - \sum_{n=1}^N \alpha_n y_n = 0 \end{array} \right.$
 \rightarrow Sustituyendo... $\mathcal{L}(w, b, \alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m$
- 2. Maximizamos $\mathcal{L}(\alpha)$, con la restricción $\sum_{n=1}^N \alpha_n y_n = 0$. (∇_b).
 $\alpha_n \geq 0 \forall n$
- + Para ello, utilizamos programación cuadrática.
- * Pasamos $\mathcal{L}(\alpha)$ a representación matricial, y minimizamos $-\mathcal{L}(\alpha)$:
 $\min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - \bar{1}^T \alpha \rightarrow \sum_{n=1}^N \alpha_n$, sujeto a restricción $y^T \alpha = 0$
 $0 < \alpha < \infty \quad \sum \alpha_n y_n$
 $\begin{bmatrix} y_1 & y_2 & x_1^T x_1 & y_1 x_1^T x_2 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$

3. Calculamos, finalmente, w y b .

$$w^* = \sum_{n=1}^N \alpha_n y_n x_n = \sum_{\text{no vector de soporte}} \alpha_n y_n x_n \rightarrow \text{Si no es vector de soporte, } \alpha = 0.$$

$$b^* = y_s - \sum y_n \alpha_n x_n^T x_s.$$

\downarrow de un punto tal que margen = 1 \rightarrow Vector de soporte.

CONDICIONES DE LAS SOLUCIONES ÓPTIMAS (KKT).

4. Son características que cumplen (u^*, α^*) .

\rightarrow Problema dual.
 \rightarrow Problema "primal"

1. Restricciones $\left\{ \begin{array}{l} \text{- Primal: } \alpha_m^T u^* \geq c_m. \\ \text{- Dual: } \alpha_m > 0. \end{array} \right.$

2. "Holgura" complementaria: $\alpha_m^* (y_m (x_m^T w^* + b^*) - 1) = 0$.

* $\text{dist}(x, w) \geq 1$.
Holgura (distancia)

\rightarrow Si es = 1 \rightarrow Parte derecha = 0. $\rightarrow \alpha_m^*$ puede ser $\neq 0$.

\rightarrow Si es $> 1 \rightarrow$ Parte derecha $> 0 \rightarrow \alpha_m^* = 0$.

* "Vector de soporte": aquellos x_n tales que $\left\{ \begin{array}{l} \text{dist}(x, w) = 1 \\ \alpha_m^* > 0. \end{array} \right.$

3. Estacionariedad de u : $\nabla_u L(u, \alpha) = 0$ si $u = u^*$ y $\alpha = \alpha^*$.
($w^* > b$)

* Estas condiciones son las que garantizan que el problema dual es equivalente al primal.

TRANSFORMACIONES NO LINEALES

+ Podemos aplicar $\Phi(x) = z$, y simplemente utilizar z en vez de x .

+ ¿Cuánto perjudica esto en términos de eficiencia?

\rightarrow En vez de realizar $x_n^T x_n \rightarrow z_n^T z_n$ (producto escalar); no mucho, si z tiene una dimensión no muy grande.

+ Si visualizamos datos y separador \rightarrow los vectores de soporte no tienen por qué estar a distancia etc. (si lo están en Z -espacio).

+ ¿Cuánto afecta dim. Z en generalización? \rightarrow Es INDEPENDIENTE...

GENERALIZACIÓN DEL MODELO

+ ¿Es tener un mayor margen realmente mejor?

* Mientras que en regularización fijamos $w^T w < C$ y minim. $E_{in} \dots$

* En SVM, fijamos $E_{in} = 0$ y minimizamos $w^T w$ (equiv. maximizar $\frac{1}{\|w\|}$).

* Cuanto mayor es el margen que exigimos \rightarrow \downarrow dicotomías son separables

$\rightarrow \downarrow$ dvc. Concretamente: $dvc(p) = \min(\lceil R_p^2 \rceil, d) + 1$.
Margen exigido

* Usando validación cruzada:

\rightarrow Si no quitamos un vector soporte, la solución no cambia: $e_n = 0$.

\rightarrow Si quitamos un vector soporte \rightarrow ~~$e_n = 0$~~ $0 \leq e_n \leq 1$.

* Por tanto, $E_{out}(SVM) \leq E_{in}(SVM) = \frac{1}{N} \sum_{n=1}^N e_n \leq \frac{\# \text{ vectores soporte}}{N}$.

\rightarrow Entonces, cuantos menos vectores soporte obtengamos, mejor será la generalización (¡independientemente de dvc!).

EL TRUQUE DEL KERNEL EN SVM DUAL

+ Estamos en el caso en que usamos $\Phi: X \rightarrow Z$ sobre x :

$$L(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m \boxed{z_n^T z_m} \rightarrow \text{Producto escalar}$$

$$g(x) = \text{signo}(w^T z + b) = \text{signo}[(\sum \alpha_n y_n \boxed{z_n^T z}) + b]$$

$w = \sum \alpha_n y_n z_n$ \rightarrow Producto escalar

* Solo utilizamos z para realizar el producto escalar y obtener n^o .

* ¿Somos capaces de calcular $z^T z$ sin calcular $\Phi(x)$?

$\rightarrow z_n^T z = K(x, x')$ \leftarrow kernel (función que calcule, en algún espacio Z , $\Phi(x)^T \Phi(x')$).

+ Ejemplo: $K(x, x') = (1 + x^T x')^a$ (kernel polinómico).

$= (1 + x_1 x'_1 + \dots + x_d x'_d)^a$ $\left\{ \begin{array}{l} \text{si } a=100 \\ d=10 \end{array} \right. \rightarrow$ expansión enorme. coeficientes etc.

$=$ producto escalar en Z de dimensión altísima. ¿Qué $\Phi(x)$ es?

\rightarrow Nos da igual, existe

TEMA 7: CLASIFICADORES LINEALES ÓPTIMOS (2)

KERNEL GAUSSIANO-RBF (RADIAL BASIS FUNCTION)

$$+ k(x, x') = e^{-\gamma \|x - x'\|^2}$$

↑
Parámetro ajustable.

* Ejemplo: $d=1, \gamma=1 \rightarrow k(x, x') = e^{-(x-x')^2}$

→ Desarrollo de Taylor: $= \exp(-x^2) \exp(-x'^2) \sum_{k=0}^{\infty} \frac{2^k (x)^k (x')^k}{k!} \exp(2xx')$

* Podemos ver que $\begin{cases} \bullet \text{ es } f(x) \\ \bullet \text{ es } f(x') \\ \bullet \text{ es cte.} \end{cases}$

* Si agrupamos, tenemos un producto escalar entre $f(x)$ y $f(x')$ de ∞ términos.

* Este kernel tiene dimensión infinita $\rightarrow \Phi(x)$ es NLT infinita.

CONSTRUCCIÓN DE KERNELS

+ Podemos formar nuevos kernels combinando varios de ellos.

+ Si queremos saber si k es un kernel válido \rightarrow Condición de Mercer.

Si $\begin{cases} * k(x, x') = k(x', x) \text{ (simetría)} \\ * k = \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{bmatrix} \text{ es positiva semidefinida. } (|K| > 0). \end{cases}$

SVM SUAVE

+ SVM, pero ahora toleramos errores. \rightarrow ¿Cuánto?

* Nuestro error ahora lo medimos como "violación del margen":

$$\xi_n \rightarrow y_n (w^T x_n + b) \geq 1 - \xi_n$$

"xi"

+ Minimizamos $\frac{1}{2} w^T w + C \sum_{n=1}^N \xi_n$, sujeto a $\begin{cases} y_n (w^T x_n + b) \geq 1 - \xi_n \\ \xi_n \geq 0 \forall n \in [1, N] \end{cases}$

"regularizador" $\begin{cases} C \rightarrow \infty \approx \text{Hard SVM} \\ C \rightarrow 0 \approx \text{No tenemos error en cuenta} \end{cases}$

* Este problema es muy similar a clasificación con regularización.

→ $w^T w$ función amor reg. \rightarrow C-Error función como guía.

Es cambio de la curva

DUAL DE SVM SUAVE

+ Usando Lagrangianos:

$$\mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} w^T w + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (y_n (w^T x_n + b) - 1 + \xi_n) - \sum_{n=1}^N \beta_n \xi_n$$

Añadimos el error

Como ξ tiene restricciones, surge.

* Minimizamos con w, b, ξ ; maximizamos con α, β .

$$\nabla_w \mathcal{L} = w - \sum_{n=1}^N \alpha_n y_n x_n = 0.$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_{n=1}^N \alpha_n y_n = 0 ; \quad \frac{\partial \mathcal{L}}{\partial \xi} = C - \alpha_n - \beta_n = 0$$

* Usando la fórmula anterior, tenemos que $(\sum \xi - \alpha \sum \xi - \beta \sum \xi) = 0$

$$\rightarrow \mathcal{L}(w, b, \xi, \alpha, \beta) = \frac{1}{2} w^T w - \sum_{n=1}^N \alpha_n (y_n (w^T x_n + b) - 1), \text{ con las mismas condiciones que con SVM duro, excepto: } 0 \leq \alpha_n \leq C$$

(es lo único nuevo).

+ Ahora, tenemos dos tipos de vectores de soporte:

→ En el margen $\begin{cases} 0 < \alpha_n^* < C \\ y_n (w^T x_n + b) = 1 \end{cases}$ (no tienen error).

→ No en el margen $\begin{cases} \alpha_n^* = C \\ y_n (w^T x_n + b) < 1 \end{cases}$ (hay error / viola el margen).

+ En cuanto a b^* , pasa algo interesante:

* Si solo hay un vector soporte en el margen $\rightarrow b^* = 1 - y_s (w^T x_s)$

* Pueden haber, si hay más de uno, un rango de soluciones para b .

CONCLUSIONES SOBRE SVM SUAVE

+ Puede ser visto como un caso especial de clasif. regularizada $\begin{cases} - \text{Esum}(b, w) \\ - \text{weight decay} \end{cases}$

+ Controla la complejidad del modelo maximizando el margen.

+ Puede lidiar con transformaciones infinitas/muy grandes \rightarrow kernel.

+ Expresa la solución con $\begin{cases} - \text{Vectores de soporte} \\ - \text{Multiplicadores de Lagrange} \\ - \text{Kernel} \end{cases}$

+ Puede controlar la sensibilidad a errores con C .

* Si C y $k(x, x')$ se eligen correctamente \rightarrow Φ y define uno de los mejores modelos de clasificación.