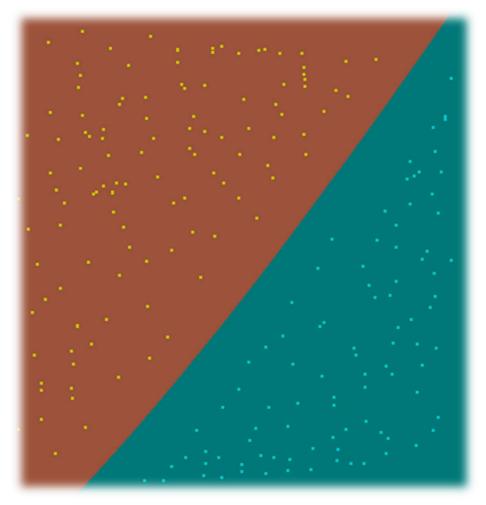
SVM, bases mathématiques

\mathbf{S} upport \mathbf{V} ector \mathbf{M} achine

Séparateur à Vaste Marge



Dans les années 1980, la star de l'IA était le langage Prolog fondé sur la logique. En ces jours, où la vague de l'IA déferle, les méthodes se fondent sur l'apprentissage statistique. Parmi ces méthodes, SVM est l'objet de ce petit document. Il existe, bien sûr de nombreuses sources sur le sujet. Mais après consultation de plusieurs d'entre elles, j'ai éprouvé une certaine frustration. Dans la plupart de ces documents, les fondements mathématiques (avec démonstrations) étaient traités assez légèrement au profit de la mise en œuvre pratique. Ce qui m'a conduit à rédiger ce texte.

Bien que ce document ne prétende nullement être un guide pratique de SVM, je fournis 3 exemples en illustration : deux exemples « jouets » à finalité pédagogique et un exemple avec des données réelles. Les outils utilisés sont le Solveur d'Excel, Tanagra et R package e1071.

- Les prérequis sont :
 - ✓ Le théorème Karush-Kuhn-Tucker version optimisation convexe.
 - ✓ Des connaissances d'algèbre linéaire, en particulier la manipulation des matrices par blocs.
 - ✓ Un minimum de familiarité avec les logiciels précités.

Optimisation quadratique

Soit h et $g_1\cdots g_n$ fonctions $\mathbb{R}^p\to\mathbb{R}$, on s'intéresse au problème d'optimisation sous contrainte (s.c.) suivant :

$$(P) \begin{cases} \min_{v \in \mathbb{R}^p} h(v) \\ s.c. \ \forall i \in 1 \dots n : g_i(v) - c_i \le 0 \end{cases}$$

En vue des applications au SVM, on supposera le problème quadratique :

- 1. $h(v) = \frac{1}{2}v^t A v + B^t v$ où $A \in \mathcal{M}_{p \times p}(\mathbb{R})$ semi-définie positive et $B \in \mathbb{R}^p$
- 2. $\forall i \in 1 \cdots n : g_i$ est une forme linéaire sur \mathbb{R}^p

Remarques:

- 1. Les conditions $g_i(v)-c_i \leq 0$ peuvent s'écrire matriciellement : $Gv-C \leq 0$ où $G \in \mathcal{M}_{n \times p}(\mathbb{R})$ et $C \in \mathbb{R}^n$.
- 2. $V = \{v \in \mathbb{R}^p \mid Gv C \le 0\}$, autrement dit l'ensemble des vecteurs qui vérifient les contraintes est un ensemble convexe non vide. Le problème (P) peut donc encore s'écrire : $\min_{v \in V} h(v)$.
- 3. $\nabla h(v) = Av + B$; $\nabla^2 h(v) = A$.
- 4. La matrice A étant semi-définie positive, la fonction h est convexe. Tout minimum local est un minimum global sur V.

DEMONSTRATION:

Soit \tilde{v} un minimum local de f sur V. Supposons qu'il existe $v \in V$ tel que $f(v) \le f(\tilde{v})$ et posons $y_t = tv + (1-t)\tilde{v}$ $t \in]01[$. Alors par convexité de f, $f(y_t) \le t f(v) + (1-t) f(\tilde{v})$. Pour $t \in]01[$ t, $f(v) + (1-t) f(\tilde{v}) < f(\tilde{v})$, d'où : $f(y_t) < f(\tilde{v})$.

Par ailleurs, il existe un voisinage U de \tilde{v} tel que pour tout $x \in U$, $f(\tilde{v}) \leq f(x)$. Pour t suffisamment petit, $y_t \in U$ et alors $f(\tilde{v}) \leq f(y_t)$, contradiction.

5. Si, de plus, la matrice A est définie positive, la fonction h est strictement convexe et la solution est unique.

DEMONSTRATION:

Supposons qu'il existe deux minimums distincts \tilde{v}_1 et \tilde{v}_2 , selon 4., ils sont globaux et $f\left(\tilde{v}_1\right) = f\left(\tilde{v}_2\right) = m$. Par convexité stricte : $f\left(\frac{\tilde{v}_1 + \tilde{v}_2}{2}\right) < \frac{f\left(\tilde{v}_1\right) + f\left(\tilde{v}_2\right)}{2} = m$, contradiction.

Karush-Kuhn-Tucker

Dans le cadre de l'optimisation avec h convexe de classe C^2 et des contraintes affines on peut appliquer le théorème Karush-Kuhn-Tucker version « forte » :

$$\tilde{v} \text{ solution de } (P)$$

$$\tilde{\varphi}$$

$$\exists \tilde{\alpha} \in \mathbb{R}^{n}_{+} \text{ tel que :}$$

$$\nabla h(\tilde{v}) + \sum_{i=1}^{n} \tilde{\alpha}_{i} \nabla g_{i}(\tilde{v}) = 0 \qquad (1)$$

$$\text{et}$$

$$\forall i \in 1 \cdots n : \tilde{\alpha}_{i} (g_{i}(\tilde{v}) - c_{i}) = 0 \qquad (2)$$

Les $\tilde{\alpha}_i$ se nomment les multiplicateurs de Lagrange. L'interprétation de cette deuxième égalité est importante. Elle signifie que si $0 < \tilde{\alpha}_i$, alors la contrainte est saturée $\left(g_i(\tilde{v}) - c_i = 0\right)$.

Lagrangien

Le lagrangien du problème est alors défini pour $v \in \mathbb{R}^p$, $\alpha \in \mathbb{R}^n_+$ par :

$$\mathcal{L}(v,\alpha) = h(v) + \sum_{i=1}^{n} \alpha_{i} (g_{i}(v) - c_{i})$$

Or, pour $\alpha \in \mathbb{R}^n_+$ et pour $v \in V$: $g_i(v) - c_i \le 0$, donc $\sum_{i=1}^n \alpha_i (g_i(v) - c_i) \le 0$. Il en résulte que :

Pour
$$v \in V$$
, $\max_{\alpha \in \mathbb{R}^n_+} \mathcal{L}(v, \alpha) = h(v)$.

Dans le cadre de l'optimisation convexe, on a :

$$\min_{v \in \mathbb{R}^{n}} \max_{\alpha \in \mathbb{R}^{n}_{+}} \mathcal{L}(v, \alpha) = \max_{\alpha \in \mathbb{R}^{n}_{+}} \min_{v \in \mathbb{R}^{p}} \mathcal{L}(v, \alpha)$$

Le problème (P), dit primal peut donc s'écrire : $\min_{v \in V} \max_{\alpha \in \mathbb{R}^n_+} \mathcal{L}(v, \alpha)$ et est équivalent au problème

$$(D), \, \mathrm{dit} \, \, \mathrm{dual} : \, \max_{lpha \in \mathbb{R}^n_+} \min_{v \in \mathbb{R}^p} \mathcal{L} ig(v, lpha ig).$$

Point-selle du Lagrangien

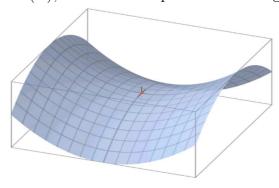
 $(\tilde{v}, \tilde{\alpha}) \in \mathbb{R}^p \times \mathbb{R}^n_+$ est un point-selle du lagrangien si et seulement si :

$$\forall v \in \mathbb{R}^p, \forall \alpha \in \mathbb{R}^n : \mathcal{L}(\tilde{v}, \alpha) \leq \mathcal{L}(\tilde{v}, \tilde{\alpha}) \leq \mathcal{L}(v, \tilde{\alpha})$$

Autrement dit:

$$\mathcal{L}(\tilde{v}, \tilde{\alpha}) = \max_{\alpha \in \mathbb{R}^{n}} \mathcal{L}(\tilde{v}, \alpha) = \min_{v \in \mathbb{R}^{p}} \mathcal{L}(v, \tilde{\alpha})$$

Si $(\tilde{v}, \tilde{\alpha})$ est solution de (P) ou (D), c'est alors un point-selle du lagrangien.



Retour au problème quadratique

$$\checkmark$$
 Si $h(v) = \frac{1}{2}v^t Av + B^t v$, alors : $\nabla h(v) = Av + B$.

$$\checkmark$$
 Si $g_i(v) = G_{i,\cdot}v$ forme linéaire, alors $\nabla g_i(v) = G_{i,\cdot}^t$ et $\sum_{i=1}^n \tilde{\alpha}_i \nabla g_i(v) = G^t \tilde{\alpha}$.

L'égalité (1) s'écrit alors matriciellement : $A\tilde{v} + B + G^t \tilde{\alpha} = 0$.

L'égalité (2) s'écrit alors matriciellement : $Diag\left(\tilde{\alpha}\right)\left(G\tilde{v}-C\right)=0$

Le lagrangien s'écrit alors :

$$\mathcal{L}(v,\alpha) = h(v) + \sum_{i=1}^{n} \alpha_i (g_i(v) - c_i) = \frac{1}{2} v^t A v + B^t v + \alpha^t (Gv - C)$$

Résumé optimisation quadratique $A \in \mathcal{M}_{p}(\mathbb{R})$ semi-définie positive, $B \in \mathbb{R}^{p}, G \in \mathcal{M}_{p \times p}(\mathbb{R})$				
Primal Dual				
$\begin{cases} \min_{v \in \mathbb{R}^p} \left(\frac{1}{2} v^t A v + B^t v \right) \\ s.c. G v - C \le 0 \end{cases}$	$\max_{\alpha \in \mathbb{R}_{+}^{n}} (\mathcal{L}_{D}(\alpha))$ Où $\mathcal{L}_{D}(\alpha) = \min_{v \in \mathbb{R}^{p}} \left(\frac{1}{2} v^{t} A v + B^{t} v + \alpha^{t} (G v - C) \right)$			
Si \tilde{v} solution du primal Si $\tilde{\alpha}$ solution du dual				
\tilde{v} et $\tilde{\alpha}$ sont liés par les relations : $A\tilde{v} + B + G^{t}\tilde{\alpha} = 0$ (3) et $Diag(\tilde{\alpha})(G\tilde{v} - C) = 0$ (4)				

Application à SVM

Deux nuages linéairement séparables, marge rigide (hard-margin)

Données:

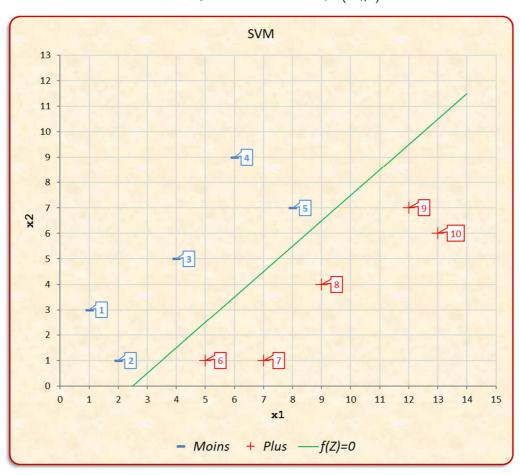
$$X \in M_{n \times p}(\mathbb{R})$$
 n vecteurs (ou points) dans $(\mathbb{R}^p)^*$ à chaque point $X_{i,\bullet} = (x_{i,1}, \dots, x_{i,p}) \in (\mathbb{R}^p)^*$ on associe une étiquette $y_i \in \{-1, 1\}$

Chacune des p colonnes de X correspond à une variable numérique, chacune des n ligne, aux valeurs des variables pour un individu.

Avertissement

On a un hiatus entre deux conventions. D'une part l'algèbre linéaire où les éléments de \mathbb{R}^p sont représentés par des vecteurs colonnes. D'autre part la statistique où dans un tableau de données $(data\ frame)$ les vecteurs décrivant les individus sont écrits en ligne. Compte tenu de l'isomorphisme canonique d'espaces euclidiens entre \mathbb{R}^p et son dual $(\mathbb{R}^p)^*$, cela n'a pas de conséquence sauf dans les écritures matricielles du genre $f(Z) = Z\beta + \beta_0 = 0$ où Z doit être un vecteur-ligne et β un vecteur colonne.

On suppose que les données sont linéairement séparables, c'est-à-dire qu'il existe un hyperplan $H_0 \text{ de } \left(\mathbb{R}^p\right)^* \text{ d'équation } f\left(Z\right) = Z\beta + \beta_0 = 0 \text{ tel que } 0 < y_i f\left(X_{i,\bullet}\right).$



Ensuite on se met en quête du « meilleur » hyperplan séparateur. En langage imagé et en dimension 2, on le définit comme celui qui va définir la « route » la plus large possible entre les deux nuages. Les marges de cette « route » sont définies par les équations f(Z) = a et f(Z) = -a de deux hyperplans H_1 et H_{-1} parallèles à H_0 . Comme les coefficients des équations sont définies à une constante multiplicative près, les équations peuvent se ramener à : f(Z) = 1 et f(Z) = -1 Les conditions de séparation s'expriment : $\forall i \in 1 \cdots n \colon 1 \leq y_i \left(X_{i,\bullet} \beta + \beta_0\right)$.

La distance d'un point M à l'hyperplan H_0 est donné par : $d(M, H_0) = \frac{|M\beta + \beta_0|}{\|\beta\|}$. Si $M \in H_1$, $M\beta + \beta_0 = 1$ et $d(M, H_0) = \frac{1}{\|\beta\|}$. La marge qui sépare les deux hyperplans H_1 et H_{-1} est donc $d(H_1, H_{-1}) = \frac{2}{\|\beta\|}$ La maximalisation de cette marge conduit à résoudre le problème d'optimisation suivant :

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \in \mathbb{R}^p, \, \beta_0 \in \mathbb{R},$$

$$\min_{\beta, \beta_0} \left(\frac{\|\beta\|^2}{2} \right) \text{ s.c. } \forall i \in 1 \dots n \colon 1 \leq y_i \left(X_{i, \bullet} \beta + \beta_0 \right)$$

Il s'agit d'un problème d'optimisation quadratique sous contrainte affine avec :

$$v = \left(\frac{\beta_0}{\beta}\right) \in \mathbb{R}^{1+p} ; A = Diag\left(0, 1 \cdots 1\right) \in \mathcal{M}_{\left(1+p\right) \times \left(1+p\right)}\left(\mathbb{R}\right) ; B = 0$$

$$G = -Diag(y)(\vec{1} \mid X) \in \mathcal{M}_{n \times (p+1)}(\mathbb{R}) \quad ; C = -\vec{1} \text{ où } \vec{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}, \ y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n$$

L'égalité (3) donne :

$$A\tilde{v} + B + G^{t}\tilde{\alpha} = 0 \Rightarrow \left(\frac{0}{\tilde{\beta}}\right) - \left(\frac{\vec{\mathbf{1}}^{t}}{X^{t}}\right) Diag\left(y_{1}, \dots, y_{n}\right) \tilde{\alpha} = 0 \Rightarrow \begin{cases} \sum_{i=1}^{n} y_{i} \tilde{\alpha}_{i} = 0 \Leftrightarrow y^{t} \tilde{\alpha} = 0 \text{ (5)} \\ \tilde{\beta} = \sum_{i=1}^{n} \tilde{\alpha}_{i} y_{i} X_{i, \cdot}^{t} \Leftrightarrow \tilde{\beta} = X^{t} Diag(y) \tilde{\alpha} \text{ (6)} \end{cases}$$

L'égalité (2) ou (4) donne :

$$\forall i \in 1 \cdots n : \tilde{\alpha}_{i} \left(y_{i} \left(X_{i, \bullet} \tilde{\beta} + \tilde{\beta}_{0} \right) - 1 \right) = 0 \Leftrightarrow Diag(\alpha) \left(\vec{1} - Diag(y) \left(\beta_{0} \vec{1} \mid X \beta \right) \right) = 0$$

Lagrangien

Le lagrangien s'écrit :

$$\mathcal{L}(v,\alpha) = \frac{1}{2}\beta^{t}\beta + \sum_{i=1}^{n}\alpha_{i}\left(-y_{i}\left(X_{i,\bullet}\beta + \beta_{0}\right) + 1\right) = \frac{1}{2}\beta^{t}\beta - \sum_{i=1}^{n}\alpha_{i}y_{i}X_{i,\bullet}\beta - \sum_{i=1}^{n}\alpha_{i}y_{i}\beta_{0} + \sum$$

Pour $(\tilde{v}, \tilde{\alpha})$ point-selle du lagrangien

$$\mathcal{L}(\tilde{v}, \tilde{\alpha}) = \frac{1}{2} \tilde{\beta}^t \tilde{\beta} - \tilde{\alpha}^t Diag(y) X \tilde{\beta} - \tilde{\beta}_0 \underbrace{\sum_{i=1}^n \tilde{\alpha}_i y_i}_{=0 (5)} + \tilde{\alpha}^t \tilde{\mathbf{1}}.$$

En utilisant (6): $\mathcal{L}(\tilde{v}, \tilde{\alpha}) = \frac{1}{2} \tilde{\alpha}^t Diag(y) XX^t Diag(y) \tilde{\alpha} - \alpha^t Diag(y) XX^t Diag(y) \tilde{\alpha} + \alpha^t \vec{1}$. Au final:

$$\mathcal{L}(\tilde{v}, \tilde{\alpha}) = \tilde{\alpha}^t \vec{1} - \frac{1}{2} \tilde{\alpha}^t Diag(y) XX^t Diag(y) \tilde{\alpha}$$

Or
$$\mathcal{L}(\tilde{v}, \alpha) \leq \mathcal{L}(\tilde{v}, \tilde{\alpha})$$
, donc $\tilde{\alpha} = \underset{\alpha \in \mathbb{R}^{n}}{\arg \max} \left(\alpha^{t} \vec{\mathbf{1}} - \frac{1}{2} \alpha^{t} Diag(y) XX^{t} Diag(y) \alpha \right)$

Le problème dual s'écrit donc :

$$\max_{\alpha \in \mathbb{R}^{n}_{+}} \left(\alpha^{t} \vec{\mathbf{1}} - \frac{1}{2} \alpha^{t} Diag(y) XX^{t} Diag(y) \alpha \right) \text{ s.c. } \alpha^{t} y = 0$$

Vecteur support

Les égalités : $\forall i \in 1 \cdots n : \tilde{\alpha}_i \left(y_i \left(X_{i,\bullet} \tilde{\beta} + \tilde{\beta}_0 \right) - 1 \right)$ montrent que si $0 < \tilde{\alpha}_i$, alors la contrainte i est saturée : $y_i \left(X_{i,\bullet} \tilde{\beta} + \tilde{\beta}_0 \right) = 1$. On définit alors $S = \left\{ i \in 1 \cdots n / 0 < \alpha_i \right\}$. Les $X_{i,\bullet}, i \in S$, sont appelés vecteurs supports.

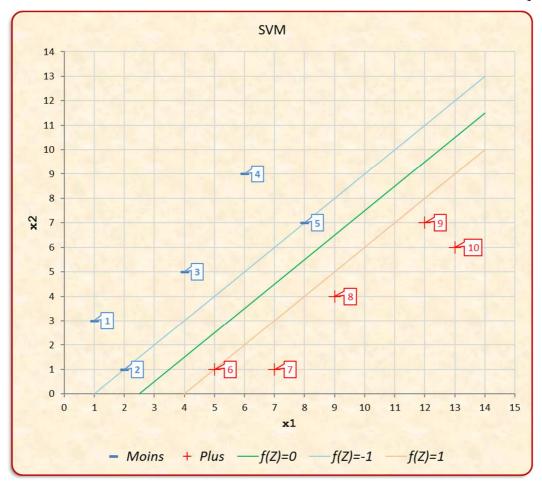
Du dual, retour vers le primal

Soit $\tilde{\alpha} \in \mathbb{R}^n_+$ solution du dual, alors (6) permet de calculer $\tilde{\beta} = \sum_{i=1}^n \tilde{\alpha}_i y_i X_{i,\cdot}^t = \sum_{i \in S} \tilde{\alpha}_i y_i X_{i,\cdot}^t$.

Si
$$i \in S$$
 $y_i \left(X_{i, \bullet} \tilde{\beta} + \tilde{\beta}_0 \right) = 1 \Leftrightarrow \tilde{\beta}_0 = \frac{1 - y_i X_{i, \bullet} \tilde{\beta}}{y_i} = \underbrace{y_i - X_{i, \bullet} \tilde{\beta}}_{\text{car } v_i^2 = 1}$

Remarques:

- \checkmark Le calcul de $\tilde{m{eta}}, \tilde{m{eta}}_0$ ne fait intervenir que les <u>vecteurs supports</u>.
- ✓ Ce calcul effectué, si $Z = (z_1, \dots, z_p) \in (\mathbb{R}^p)^*$ représente les valeurs des variables pour un nouvel individu, son classement se fait à partir du signe de la valeur de la <u>fonction score</u> $f(Z) = Z\tilde{\beta} + \tilde{\beta}_0$.



Les vecteurs supports sont 6, 5 et 2

Résumé marge rigide					
Primal	Dual				
$\beta = \begin{pmatrix} \beta_{1} \\ \vdots \\ \beta_{p} \end{pmatrix} \in \mathbb{R}^{p}, \beta_{0} \in \mathbb{R}, f(X_{i,\bullet}) = X_{i,\bullet}\beta + \beta_{0}$ $\min_{\substack{\beta \in \mathbb{R}^{p} \\ \beta_{0} \in \mathbb{R}}} \left(\frac{\ \beta\ ^{2}}{2} \right) \text{ sous la contrainte} : y_{i}f(X_{i,\bullet}) \ge 1$	$\max_{\alpha \in \mathbb{R}^{n}_{+}} \mathcal{L}_{D}(\alpha) \text{ s.c.} : \sum_{i=1}^{n} \alpha_{i} y_{i} = 0$ $\mathcal{L}_{D}(\alpha) = \alpha^{t} \vec{1} - \frac{1}{2} \alpha^{t} Diag(y) X X^{t} Diag(y) \alpha$ $\mathcal{L}_{D}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{i=1}^{n} \alpha_{i} y_{i} \langle X_{i,\bullet}; X_{j,\bullet} \rangle y_{j} \alpha_{j}$ $S \text{ ensemble des } i \in 1 \cdots n \ \alpha_{i} \neq 0$ $\text{Si } i \in S, X_{i,\bullet} \text{ est un vecteur support}$ $\tilde{\beta} = \sum_{i=1}^{n} \tilde{\alpha}_{i} y_{i} X_{i,\bullet}^{t} \text{ et } \tilde{\beta}_{0} = y_{i} - X_{i,\bullet} \tilde{\beta} \ (i \in S)$				
$p - \sum_{i \in S} \alpha_i y_i X_{i,\bullet} \text{et} p_0 - y_i - X_{i,\bullet} p (i \in S)$					
Fonction score: $f(Z) = Z\tilde{\beta} + \tilde{\beta}_0 = \sum_{i \in S} \tilde{\alpha}_i y_i \langle X_{i, \bullet}; Z \rangle + \tilde{\beta}_0$					

Deux nuages non linéairement séparables,

Avec des vraies données, la séparabilité linéaire est peu probable. Pour y remédier, on introduit des variables d'écart ξ_i (slack variables ou variables ressort) et la condition de « bon classement » $1 \le y_i \left(X_{i,\bullet} \beta + \beta_0 \right)$ est modifiée en $1 - \xi_i \le y_i \left(X_{i,\bullet} \beta + \beta_0 \right)$ afin de tolérer le mauvais classement de certains points :

$$\begin{aligned} &\xi_i = 0 \text{ point } i \text{ du bon côt\'e de sa marge} \\ &0 < \xi_i < 1 \text{ point } i \text{ entre les 2 marges} \\ &1 \le \xi_i \text{ point } i \text{ mal class\'e} \end{aligned}$$

Le problème s'énonce alors :

$$\beta = \begin{pmatrix} \beta_{1} \\ \vdots \\ \beta_{p} \end{pmatrix} \in \mathbb{R}^{p}, \ \beta_{0} \in \mathbb{R}, \xi = \begin{pmatrix} \xi_{1} \\ \vdots \\ \xi_{n} \end{pmatrix} \in \mathbb{R}^{n}$$

$$\min_{\beta, \beta_{0}, \xi} \left(\frac{\|\beta\|^{2}}{2} + \Gamma \sum_{i=1}^{n} \xi_{i} \right) \text{ s.c. } \forall i \in 1 \dots n \colon 1 - \xi_{i} \leq y_{i} \left(X_{i, \bullet} \beta + \beta_{0} \right) \text{ et } 0 \leq \xi_{i}$$

 $0 \le \Gamma^1$ est un paramètre dit de coût (cost) qui règle la tolérance aux erreurs.

Là encore il s'agit d'un problème d'optimisation quadratique avec :

Rappel de notations :
$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \vec{1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^n$$

$$v = \left(\frac{\beta_{0}}{\beta}\right) \in \mathbb{R}^{1+p+n}; A = Diag\left(0, \underbrace{1 \cdots 1}_{p \text{ fois}}, \underbrace{0 \cdots 0}_{n \text{ fois}}\right) \in \mathcal{M}_{(1+p+n) \times (1+p+n)}(\mathbb{R}); B = \Gamma\begin{pmatrix} 0 \\ \vdots \\ 0 \\ \overline{1} \end{pmatrix} \in \mathbb{R}^{1+p+n}$$

Les n premières contraintes s'écrivent :

$$-y_i\beta_0 - y_iX_{i,\bullet}\beta - \xi_i \le -1$$

Les n suivantes :

$$-\xi_i \leq 0$$

Ce qui se traduit par :

$$G = -\left(\frac{y \mid Diag(y)X \mid I}{0 \mid 0 \mid I}\right) \in \mathcal{M}_{2n \times (n+p+1)}(\mathbb{R}) ; C = \begin{pmatrix} -\vec{1} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{2n}$$

Comme on a 2n contraintes, on aura 2n multiplicateurs de Lagrange que l'on notera : $\begin{pmatrix} \tilde{\alpha} \\ \bar{\delta} \end{pmatrix} \in \mathbb{R}^{2n}_+$

 $^{^{1}}$ Ce paramètre est traditionnellement noté $\boldsymbol{\mathcal{C}}$, mais il est déjà utilisé dans le texte.

Si
$$\tilde{v} = \begin{pmatrix} \frac{\tilde{\beta}_0}{\tilde{\beta}} \\ \frac{\tilde{\xi}}{\tilde{\xi}} \end{pmatrix}$$
 solution du primal, et $\begin{pmatrix} \frac{\tilde{\alpha}}{\tilde{\delta}} \end{pmatrix} = \in \mathbb{R}^{2n}_+$ solution du dual, l'égalité (3) donne :

$$A\tilde{v} + B + G^{t} \begin{pmatrix} \tilde{\alpha} \\ \bar{\delta} \end{pmatrix} = 0 \Rightarrow \begin{pmatrix} 0 \\ \frac{\bar{\beta}}{\bar{\mathbf{0}}} \end{pmatrix} + \Gamma \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ \bar{\mathbf{1}} \end{pmatrix} - \begin{pmatrix} y^{t} & 0 \\ \overline{X}^{t} Diag(y) & 0 \\ \overline{I} & I \end{pmatrix} \begin{pmatrix} \tilde{\alpha} \\ \bar{\delta} \end{pmatrix} = 0 \Rightarrow \begin{pmatrix} y^{t} \tilde{\alpha} = 0 \Leftrightarrow \sum_{i=1}^{n} y_{i} \tilde{\alpha}_{i} = 0 \end{cases} \Rightarrow \begin{pmatrix} \tilde{\beta} = X^{t} Diag(y) \tilde{\alpha} \Leftrightarrow \tilde{\beta} = \sum_{i=1}^{n} \tilde{\alpha}_{i} y_{i} X_{i, \cdot}^{t} \end{cases}$$

$$\Gamma \tilde{\mathbf{1}} = \tilde{\alpha} + \tilde{\delta} \Leftrightarrow \forall i \in 1 \cdots n : \tilde{\delta}_{i} = \Gamma - \tilde{\alpha}_{i} \end{cases}$$

$$(5)$$

REMARQUES:

- \checkmark Il n'est nul besoin de connaitre $\tilde{\delta}$ pour calculer $\tilde{\beta}$.
- ✓ L'égalité (7) impose $\tilde{\alpha} \leq \Gamma \vec{1}$.

En appliquant (4) $Diag(\tilde{\alpha})(G\tilde{v}-C)=0$, on obtient:

$$Diag\begin{pmatrix} \tilde{\alpha} \\ \tilde{\delta} \end{pmatrix} \left(-\begin{pmatrix} y & Diag(y)X & I \\ 0 & 0 & I \end{pmatrix} \begin{pmatrix} \tilde{\beta}_{0} \\ \frac{\tilde{\beta}}{\tilde{\xi}} \end{pmatrix} + \begin{pmatrix} \tilde{\mathbf{1}} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \right) = 0 \Leftrightarrow \begin{cases} Diag(\tilde{\alpha}) \left(-y\tilde{\beta}_{0} - Diag(y)X\tilde{\beta} - \tilde{\xi} + \tilde{\mathbf{1}} \right) = 0 \\ Diag(\tilde{\delta})\tilde{\xi} = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} \forall i \in 1 \cdots n : \tilde{\alpha}_{i} \left(y_{i} \left(X_{i,.}\tilde{\beta} + \tilde{\beta}_{0} \right) - \left(1 - \tilde{\xi}_{i} \right) \right) = 0 \text{ (8)} \\ \forall i \in 1 \cdots n : \tilde{\xi}_{i} \left(\Gamma - \tilde{\alpha}_{i} \right) = 0 \text{ (9)} \end{cases}$$

Lagrangien

$$\mathcal{L}(v,\alpha,\delta) = \frac{1}{2}\beta^{i}\beta + \Gamma\sum_{i=1}^{n}\xi_{i} - \sum_{i=1}^{n}\alpha_{i}\left(y_{i}\left(X_{i,\bullet}\beta + \beta_{0}\right) - 1 + \xi_{i}\right) - \sum_{i=1}^{n}\delta_{i}\xi_{i}$$

$$= \frac{1}{2}\beta^{i}\beta + \Gamma\sum_{i=1}^{n}\xi_{i} - \sum_{i=1}^{n}\alpha_{i}y_{i}\left(X_{i,\bullet}\beta\right) - \beta_{0}\sum_{i=1}^{n}\alpha_{i}y_{i} + \sum_{i=1}^{n}\alpha_{i} - \sum_{i=1}^{n}\left(\alpha_{i} + \delta_{i}\right)\xi_{i}$$

$$= \frac{1}{2}\beta^{i}\beta + \Gamma\sum_{i=1}^{n}\xi_{i} - \alpha^{i}Diag\left(y\right)X\beta - \beta_{0}\sum_{i=1}^{n}\alpha_{i}y_{i} + \sum_{i=1}^{n}\alpha_{i} - \sum_{i=1}^{n}\left(\alpha_{i} + \delta_{i}\right)\xi_{i}$$

Pour $(\tilde{v}, \tilde{\alpha}, \tilde{\delta})$ point-selle du lagrangien, en tenant compte de (7) et (5):

$$\mathcal{L}(\tilde{v}, \tilde{\alpha}) = \frac{1}{2} \tilde{\beta}^{t} \tilde{\beta} + \Gamma \sum_{i=1}^{n} \tilde{\xi}_{i} - \tilde{\alpha}^{t} Diag(y) X \tilde{\beta} - \tilde{\beta}_{0} \sum_{i=1}^{n} \tilde{\alpha}_{i} y_{i} + \sum_{i=1}^{n} \tilde{\alpha}_{i} - \sum_{i=1}^{n} (\tilde{\alpha}_{i} + \tilde{\delta}_{i}) \tilde{\xi}_{i}$$

$$= \tilde{\alpha}^{t} \vec{1} + \frac{1}{2} \tilde{\beta}^{t} \tilde{\beta} - \tilde{\alpha}^{t} Diag(y) X \tilde{\beta}$$

Dès lors il suffit d'utiliser le (6) pour conclure comme dans le cas rigide. Le problème dual s'écrit alors :

$$\max_{\alpha \in \mathbb{R}^{n}_{+}} \left(\alpha^{t} \vec{\mathbf{1}} - \frac{1}{2} \alpha^{t} Diag(y) XX^{t} Diag(y) \alpha \right) \text{ s.c. } \alpha^{t} y = 0 \text{ et } \alpha \leq \Gamma \vec{\mathbf{1}}$$

Vecteur support

Comme pour le cas rigide, l'égalité (8) permet de définir les vecteurs supports par $0 < \tilde{\alpha}_i$, ce qui implique la saturation de la contrainte : $y_i \left(X_{i,\cdot} \tilde{\beta} + \tilde{\beta}_0 \right) = 1 - \tilde{\xi}_i$. Les vecteurs supports se partagent en deux types :

- \checkmark Type S_I : $\tilde{\alpha}_i < \Gamma$, ce qui implique, selon (9) la nullité de la variable d'écart $\tilde{\xi}_i = 0 \Rightarrow y_i \left(X_{i, \bullet} \tilde{\beta} + \tilde{\beta}_0 \right) = 1$. Le vecteur se situe sur une marge.
- \checkmark Type $S_{II}: \tilde{\alpha}_i = \Gamma$.

Du dual, retour vers le primal

Soit $\tilde{\alpha} \in \mathbb{R}^n_+$ solution du dual, alors (6) permet de calculer $\tilde{\beta} = \sum_{i=1}^n \tilde{\alpha}_i y_i X_{i,\cdot}^t = \sum_{i \in S} \tilde{\alpha}_i y_i X_{i,\cdot}^t$.

Pour calculer $\tilde{\beta}_0$, il faut utiliser un vecteur support de type I qui vérifie $y_i(X_i, \tilde{\beta} + \tilde{\beta}_0) = 1$.

Pour $Z \in (\mathbb{R}^p)^*$, on a alors la fonction score : $f(Z) = Z\tilde{\beta} + \tilde{\beta}_0$.

Les $\tilde{\xi}_i$ peuvent se calculer ainsi :

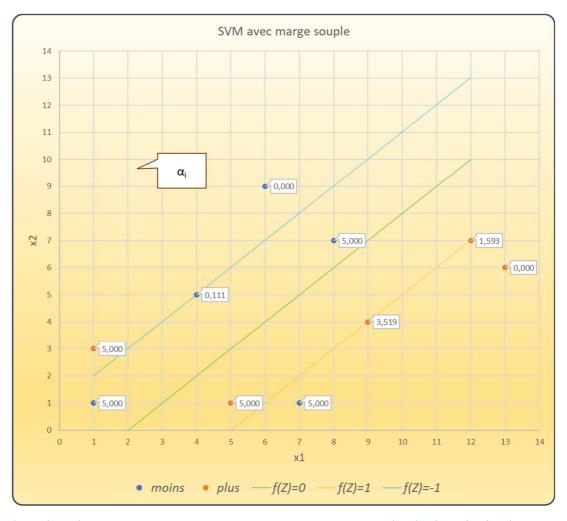
- 1. Si $X_{i,\bullet}$ est un vecteur support $\tilde{\xi}_i = 1 y_i (X_{i,\bullet} \tilde{\beta} + \tilde{\beta}_0)$.
- 2. Sinon $\tilde{\alpha}_i = 0$ et (9) implique $\tilde{\xi}_i = 0$.

On a donc : $\xi_i = \max\left(0, 1 - y_i\left(X_{i,\bullet}\tilde{\beta} + \tilde{\beta}_0\right)\right)$

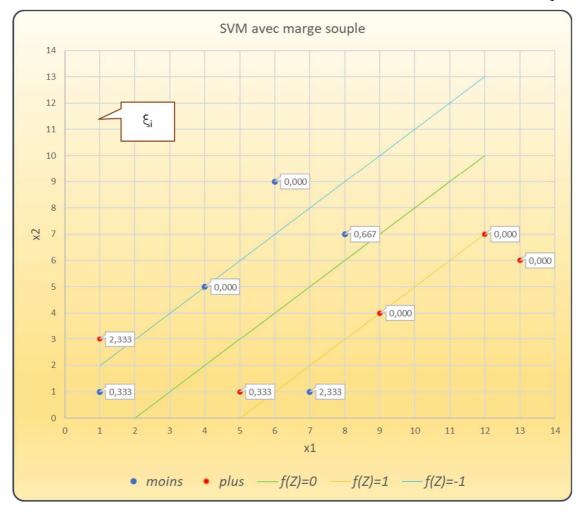
Résumé marge souple				
Primal	Dual			
$ \min_{\substack{\beta \in \mathbb{R}^{p}, \\ \beta_{0} \in \mathbb{R}, \\ \xi_{i} \in \mathbb{R}^{n}}} \left(\frac{\ \beta\ ^{2}}{2} + \Gamma \sum_{i=1}^{n} \xi_{i} \right) \Gamma \text{ constante à fixer } $ $ \text{s.c. } : 1 - \xi_{i} \leq y_{i} \left(X_{i, \bullet} \beta + \beta_{0} \right) \text{ et } 0 \leq \xi_{i} $ $ \xi_{i} : \text{ variables d'écart } $ $ \xi_{i} = 0 \text{ point } i \text{ du bon côté de sa marge } $ $ 0 < \xi_{i} < 1 \text{ point } i \text{ entre les 2 marges } $ $ 1 \leq \xi_{i} \text{ point } i \text{ mal classé } $ $ \xi_{i} = \max \left(0, 1 - y_{i} \left(X_{i, \bullet} \tilde{\beta} + \tilde{\beta}_{0} \right) \right) $	$\max_{\alpha \in \mathbb{R}^{n}_{+}} \mathcal{L}_{D}(\alpha) \text{ s.c.} : 0 \leq \alpha_{i} \leq \Gamma \text{ et } \sum_{i=1}^{n} \alpha_{i} y_{i} = 0$ $\mathcal{L}_{D}(\alpha) = \alpha^{i} \vec{1} - \frac{1}{2} \alpha^{i} Diag(y) X X^{i} Diag(y) \alpha$ $\mathcal{L}_{D}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} y_{i} \left\langle X_{i,\bullet}; X_{j,\bullet} \right\rangle y_{j} \alpha_{j}$ $S \text{ ensemble des } i \in 1 \cdots n \ \alpha_{i} \neq 0$ $\text{Si } i \in S, X_{i,\bullet} \text{ est un vecteur support}$ $\text{de type I si } \xi_{i} = 0, \text{ de type II sinon}$ $\tilde{\beta} = \sum_{i \in S} \alpha_{i} y_{i} X_{i,\bullet}^{t} \text{ et } \tilde{\beta}_{0} = y_{i} - X_{i,\bullet} \tilde{\beta} \ (i \in S_{I})$			
Fonction score: $f(Z) = Z\tilde{\beta} + \tilde{\beta}_0 = \sum_{i \in S} \tilde{\alpha}_i y_i \langle X_{i, \bullet}; Z \rangle + \tilde{\beta}_0$				

 $Exemple\ d'après\ Ricco\ Rakotomalala\ (classeur\ Excel\ ex1_dual.xlsx)$

individus	X ₁	X ₂	у	$ ilde{lpha}$	Zvs	$f(x_1, x_2)$	Γ	= 5
1	1	1	-1	5,000	0,333	-0,667		
2	4	5	-1	0,111	0,000	-1,000		S_I
3	6	9	-1	0,000	0,000	-1,667		
4	8	7	-1	5,000	0,667	-0,333	,	SII
5	7	1	-1	5,000	2,333	1,333		
6	1	3	1	5,000	2,333	-1,333	$ ilde{oldsymbol{eta}}_{\!\scriptscriptstyle 0}$	-0,667
7	5	1	1	5,000	0,333	0,667	õ	0,333
8	13	6	1	0,000	0,000	1,667	$ ilde{oldsymbol{eta}}$	-0,333
9	9	4	1	3,519	0,000	1,000		
10	12	7	1	1,593	0,000	1,000		



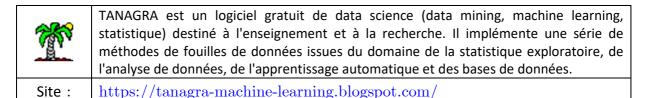
Sauf (6,9) et (13,6) tous les vecteurs sont supports, mais seuls (4,5), (12,7), (9,4) sont de type I



 $\xi_i = 0$ point i du bon côté de sa marge $0 < \xi_i < 1$ point i entre les 2 marges $1 \le \xi_i$ point i mal classé

Au sujet du traitement de cet exemple

Les calculs ont été effectués d'une part avec solveur d'Excel², de l'autre par le logiciel Tanagra de Ricco Rakotomalala.



Un des avantages pratiques de Tanagra est son intégration comme complément d'Excel (Tanagra.xla) Après sélection des données, il suffit d'aller dans le menu **Compléments** d'Excel pour lancer Tanagra. Autre avantage, les résultats sont retournés sous forme de page HTML, que l'on peut par un copier/coller intégrer directement dans une feuille Excel ou un document Word. C'est ce qui a été fait pour la page suivante.

² Cf. ex1_dual.xlsx

Supervised Learning 1 (SVM)

Parameters SVM Parameters Exponent 1 Filter type NONE Use polynom space normalization 0 Use RBF kernel 0 Gamma for RBF kernel 0,01 Complexity 5 **Calculation parameter** Epsilon for rounding 1,00E-12 Tolerance for accuracy 1,00E-03

Results

Classifier performances

Error rate				0,2		
Values prediction				Confu	sion matrix	
Value	Recall	1-Precision		moins	plus	Sum
moins	0,8	0,2	moins	4	1	5
plus	0,8	0,2	plus	1	4	5
			Sum	5	5	10

Classifier characteristics

Data description

Target attribute	y (2 values)
# descriptors	2

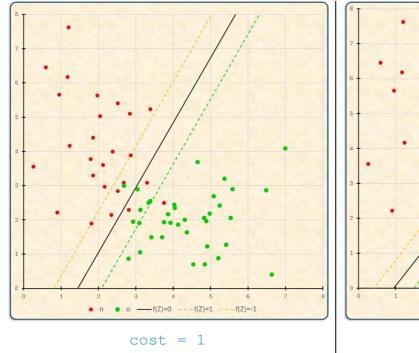
Linear classifier

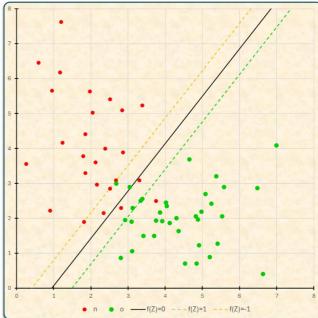
"Reference" class value : plus

	Attribute	Weight	
x1		0,333163	õ
x2		-0,332652	ρ
constant		-0,668626	$ ilde{oldsymbol{eta}}_{\!\scriptscriptstyle 0}$

Du choix de la constante Γ (cost)

Avec un petit jeu de données fabriqué pour l'occasion³, voici ce que l'on peut observer :





$$f(Z) \approx 1,57 z_1 - 0,83 z_2 - 2,27$$

$$marge = \frac{2}{\|\beta\|} \approx 1,13 \; ; \; \sum_{i=1}^{n} \xi_i \approx 9,56$$

1 • en zone rouge, 3 • en zone verte

cost = 100

$$f(Z) \approx 1.85 z_1 - 1.36 z_2 - 1.76$$

 $marge = \frac{2}{\|\beta\|} \approx 0.87 \; ; \; \sum_{i=1}^{n} \xi_i \approx 8.96$

2 • en zone rouge, 3 • en zone verte

Soit, en passant de 1 à 100 :

- ✓ Un changement dans la fonction score.
- ✓ Une réduction de la marge.
- ✓ Une réduction de la somme des écarts (variables ressorts).
- ✓ Une augmentation des mal classés.

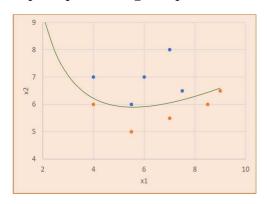
On peut avancer les explications suivantes :

La fonction à minimiser $\frac{\|\boldsymbol{\beta}\|^2}{2} + \Gamma \sum_{i=1}^n \xi_i$ est la somme de deux termes positifs. Le premier conditionne la marge, le second les écarts. Augmenter Γ , c'est mettre l'accent sur la minimisation de la somme des écarts au détriment de la minimisation de $\|\boldsymbol{\beta}\|^2$, donc de la maximalisation de la marge.

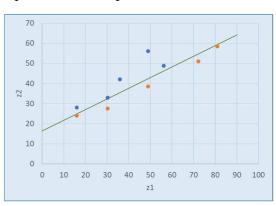
³ Cf. ex2_effet_cost.xlsx

Séparation non linéaire

La première idée est d'utiliser une transformation $\mathbb{R}^p \xrightarrow{\Phi} H$ où H désigne un espace Pré-Hilbertien réel de dimension supérieure à p (voire infini) nommé <u>espace de redescription</u>, avec l'espoir que le nuage de points transformé soit linéairement séparable. Exemple :



$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} x_1^2 \\ x_1 x_2 \end{pmatrix}$$



Une question demeure : comment déterminer Φ ?

La fonction à maximiser dans la formulation duale peut s'écrire :

$$L_{D}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{j} \left\langle X_{i,\bullet}; X_{j,\bullet} \right\rangle y_{i} y_{j}$$

Si on effectue la transformation Φ , elle s'écrira :

$$L_{D}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{j} \left\langle \Phi(X_{i,\bullet}); \Phi(X_{j,\bullet}) \right\rangle y_{i} y_{j}$$

On définit la fonction « noyau » :

$$\mathbb{R}^{p} \times \mathbb{R}^{p} \xrightarrow{K} \mathbb{R} \text{ par } K(U,V) = \langle \Phi(U); \Phi(V) \rangle = \Phi(U) \Phi(V)^{t}$$

Dans le petit exemple introductif, le noyau est : $K(U,V) = u_1 v_1 (u_1 v_1 + u_2 v_2)$

La deuxième idée est d'utiliser des fonctions noyaux K pour remplacer le produit scalaire, sans expliciter Φ .

L'astuce du noyau (kernel trick) en détail

Le problème primal transporté dans l'espace de redescription H s'écrit :

$$\min_{\substack{\beta \in H, \\ \beta_0 \in \mathbb{R}^n}} \left(\frac{\left\| \boldsymbol{\beta} \right\|^2}{2} + \Gamma \sum_{i=1}^n \xi_i \right) \text{ s.c. } : 1 - \xi_i \leq y_i \left(\left\langle \Phi \left(\boldsymbol{X}_{i, \bullet} \right) ; \boldsymbol{\beta} \right\rangle + \beta_0 \right) \text{ et } 0 \leq \xi_i$$

Et le problème dual :

$$\max_{\boldsymbol{\alpha} \in \mathbb{R}^{n}_{+}} \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} y_{i} \underbrace{\left\langle \Phi\left(X_{i,\bullet}\right); \Phi\left(X_{j,\bullet}\right) \right\rangle}_{=K\left(X_{i,\bullet}; X_{j,\bullet}\right)} y_{j} \alpha_{j} \text{ s.c. } : 0 \leq \alpha_{i} \leq \Gamma \text{ et } \sum_{i=1}^{n} \alpha_{i} y_{i} = 0$$

Soit $\tilde{\alpha}$ solution du dual, alors :

$$\begin{split} \widetilde{\boldsymbol{\beta}} &= \sum_{i \in S} \widetilde{\boldsymbol{\alpha}}_{i} y_{i} \Phi \left(\boldsymbol{X}_{i, \bullet}^{t} \right) \text{ et } \widetilde{\boldsymbol{\beta}}_{0} = y_{k} - \left\langle \Phi \left(\boldsymbol{X}_{k, \bullet}^{t} \right) ; \widetilde{\boldsymbol{\beta}} \right\rangle \; (k \in S_{\mathrm{I}}) \\ \widetilde{\boldsymbol{\beta}}_{0} &= y_{k} - \left\langle \Phi \left(\boldsymbol{X}_{k, \bullet}^{t} \right) ; \sum_{i \in S} \widetilde{\boldsymbol{\alpha}}_{i} y_{i} \Phi \left(\boldsymbol{X}_{i, \bullet}^{t} \right) \right\rangle = y_{k} - \sum_{i \in S} \widetilde{\boldsymbol{\alpha}}_{i} y_{i} \left\langle \Phi \left(\boldsymbol{X}_{k, \bullet}^{t} \right) ; \Phi \left(\boldsymbol{X}_{i, \bullet}^{t} \right) \right\rangle \\ \widetilde{\boldsymbol{\beta}}_{0} &= y_{k} - \sum_{i \in S} \widetilde{\boldsymbol{\alpha}}_{i} y_{i} K \left(\boldsymbol{X}_{k, \bullet}, \boldsymbol{X}_{i, \bullet}^{t} \right) \; (k \in S_{\mathrm{I}}) \end{split}$$

Pour $\tilde{\beta}_0$, on dispose donc d'une formule ne faisant intervenir que le noyau. En revanche, le calcul de $\tilde{\beta}$, fait intervenir la transformation Φ . Mais $\tilde{\beta}$ sert essentiellement à déterminer la fonction score f qui s'exprime ici pour $Z \in \mathbb{R}^p$: $f(Z) = \langle \Phi(Z); \tilde{\beta} \rangle + \tilde{\beta}_0$. En réinjectant l'expression de $\tilde{\beta}$ dans cette formule, on obtient :

$$f(Z) = \left\langle \Phi(Z); \sum_{i \in S} \tilde{\alpha}_i y_i \Phi(X_{i,\bullet}^t) \right\rangle + \tilde{\beta}_0 = \sum_{i \in S} \tilde{\alpha}_i y_i \left\langle \Phi(Z); \Phi(X_{i,\bullet}^t) \right\rangle + \tilde{\beta}_0 = \sum_{i \in S} \tilde{\alpha}_i y_i K(Z, X_{i,\bullet}^t) + \tilde{\beta}_0$$

Résumé dual avec marge souple et noyau

$$\max_{\alpha \in \mathbb{R}_{+}^{n}} L_{D}(\alpha) \text{ s.c.} : 0 \le \alpha_{i} \le \Gamma \text{ et } \sum_{i=1}^{n} \alpha_{i} y_{i} = 0$$

$$L_{D}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} y_{i} K(X_{i,\bullet}, X_{j,\bullet}) y_{j} \alpha_{j}$$

$$L_{D}(\alpha) = \alpha^{t} \vec{1} - \frac{1}{2} \alpha^{t} Diag(y) G(X_{1,\bullet}, \dots, X_{n,\bullet}) Diag(y) \alpha$$

S ensemble des $i\!\in\!1\!\cdots\!n$ $\tilde{\alpha}_{\scriptscriptstyle \rm i}\neq0$

Si $i \in S, X_{i,\bullet}$ est un vecteur support

de type I si $\tilde{\xi}_i = 0$, de type II sinon

$$\tilde{\beta}_0 = y_k - \sum_{i \in S} \tilde{\alpha}_i y_i K\left(X_{k,\bullet}^t, X_{i,\bullet}^t\right) \ (k \in S_{\mathrm{I}})$$

Fonction score $f(Z) = \sum_{i \in S} \tilde{\alpha}_i y_i K(Z, X_{i,\bullet}) + \tilde{\beta}_0$

On a alors : $\xi_i = \max(0, 1 - y_i f(X_{i,\bullet}))$

Quelles sont les propriétés requises pour un noyau?

On dit qu'une application symétrique $\mathbb{R}^p \times \mathbb{R}^p \xrightarrow{K} \mathbb{R}$ représente un produit scalaire dans un espace pré-hilbertien réel H de redescription si et seulement si, il existe une application $\Phi: \mathbb{R}^p \xrightarrow{\Phi} H$ telle que :

$$\forall U, V \in \mathbb{R}^p : K(U, V) = \langle \Phi(U); \Phi(V) \rangle_U$$

La réponse à la question à la question introductive tient dans la proposition suivante :

Proposition

 $\label{eq:continuous_product} \textit{Une application symétrique} \ \ \mathbb{R}^{p} \times \mathbb{R}^{p} \xrightarrow{K} \mathbb{R} \ \ \textit{représente un produit scalaire si et }$ seulement $si: \forall (X_1,\dots,X_n) \in \mathbb{R}^p \times \dots \times \mathbb{R}^p: la \ matrice \ symétrique:$

$$G\left(X_{1}, \dots, X_{n}\right) = \left(K\left(X_{i}; X_{j}\right)\right)_{i=1\dots n, j=1\dots n}$$

est semi-définie positive.

REMARQUE PREALABLE

Dans un texte sur le sujet, j'ai lu comme condition $G\left(X_{1},\cdots,X_{n}\right)$ définie positive. C'est faux et cela même dans le cas d'un noyau linéaire. Si p < n, la matrice X est de rang < n et la matrice $G(X_1,\dots,X_n) = XX^t$ aussi, donc non inversible.

DEMONSTRATION

 $1 \Rightarrow 2$

Soit K représentant un produit scalaire dans H.

$$G(X_1, \dots, X_n)$$
 semi-définie positive $\Leftrightarrow \forall U = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} : U^t G U \ge 0$. Or :

$$U^{t}GU = \sum_{i=1}^{n} \sum_{j=1}^{n} u_{i}K(X_{i,\bullet}, X_{j,\bullet})u_{j} = \sum_{i=1}^{n} \sum_{j=1}^{n} u_{i} \langle \Phi(X_{i,\bullet}); \Phi(X_{j,\bullet}) \rangle_{H} u_{j}$$

(puisque K représente un produit scalaire)

$$=\sum_{i=1}^{n}\sum_{j=1}^{n}\left\langle u_{i}\Phi\left(X_{i,\bullet}\right);u_{j}\Phi\left(X_{j,\bullet}\right)\right\rangle_{H}=\left\langle \sum_{i=1}^{n}u_{i}\Phi\left(X_{i,\bullet}\right);\sum_{j=1}^{n}u_{j}\Phi\left(X_{j,\bullet}\right)\right\rangle_{H}=\left\|\sum_{i=1}^{n}u_{i}\Phi\left(X_{i,\bullet}\right)\right\|^{2}\geq0$$

 $2 \Rightarrow 1$

Soit K un noyau vérifiant $G(X_1,\dots,X_n)$ semi-définie positive $\forall (X_1,\dots,X_n) \in \mathbb{R}^p \times \dots \times \mathbb{R}^p$.

Pour $U \in \mathbb{R}^p$, on note K_U l'application $\mathbb{R}^p \xrightarrow{K_U} \mathbb{R}$. Elle appartient à l'espace $Z \mid \longrightarrow K_U(Z) = K(U,Z)$

vectoriel des applications de $\mathbb{R}^p \to \mathbb{R}$. Soit $H = Vec(K_U/U \in \mathbb{R}^p)$. Autrement dit, H (de dimension infinie) est l'ensemble de toutes combinaisons linéaires finies des K_U .

Soit
$$\alpha = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} \in \mathbb{R}^n$$
, $\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix} \in \mathbb{R}^m$ et $J_{\alpha} = \sum_{i=1}^n \alpha_i K_{X_i}$, $J_{\beta} = \sum_{j=1}^m \beta_j K_{Y_j} \in H$, on définit :
$$\left\langle J_{\alpha} ; J_{\beta} \right\rangle_H \stackrel{def}{=} \sum_{i=1}^n \sum_{j=1}^m \alpha_i K(X_i, Y_j) \beta_j$$

Il est facile de vérifier qu'il s'agit d'une forme bilinéaire symétrique sur H.

$$\left\langle J_{\alpha};J_{\alpha}\right\rangle _{H}=\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_{i}\,K\left(X_{i},X_{j}\right)\!\alpha_{j}=\alpha^{i}G\left(X_{1},\cdots,X_{n}\right)\alpha\geq0\qquad\text{puisque}\qquad G\left(X_{1},\cdots,X_{n}\right)\qquad\text{semi-définie}$$
 positive.

 $\langle \; ; \rangle_H$ est donc une forme semi-définie positive sur H.

Pour que $\langle \; ; \; \rangle_{\!\! H} \;$ soit définie positive, autrement dit soit un produit scalaire, il reste à prouver :

$$\langle J_{\alpha}; J_{\alpha} \rangle = 0 \Rightarrow J_{\alpha} = 0$$

$$\forall J_{\alpha} \in H, U \in \mathbb{R}^{p} : \left\langle J_{\alpha}; K_{U} \right\rangle_{H} = \sum_{i=1}^{n} \alpha_{i} K(U, X_{i}) = \sum_{i=1}^{n} \alpha_{i} K_{X_{i}}(U) = J_{\alpha}(U) \quad (1)$$

Par ailleurs, l'inégalité de Cauchy-Schwarz, est vraie pour les formes semi-définies positives. D'où :

$$\forall J_{\alpha} \in H, U \in R^{p} : \left\langle J_{\alpha} ; K_{U} \right\rangle_{H}^{2} \leq \left\langle J_{\alpha} ; J_{\alpha} \right\rangle_{H} \left\langle K_{U} ; K_{U} \right\rangle_{H}.$$

Donc si $\langle J_{\alpha}; J_{\alpha} \rangle_H = 0$, alors $\langle J_{\alpha}; K(U, \bullet) \rangle_H = 0$ et alors d'après (1) $J_{\alpha}(U) = 0$ et ceci quel que soit $U \in \mathbb{R}^p$.

Pour terminer, en posant $\forall U \in \mathbb{R}^p : \Phi(U) = K_U$, on a d'après (1) :

$$\langle \Phi(U); \Phi(V) \rangle_H = \langle K_U; K_V \rangle_H = J_\alpha(V) = K(U, V)$$

COMMENTAIRE

Cette fonction Φ qui envoie dans un espace de dimension infinie, n'a guère d'intérêt pour les calculs pratiques pour lesquels l'usage du noyau qui opère en dimension fini, s'avère beaucoup plus efficace. En revanche, la proposition fournit un critère pour qu'une fonction symétrique $\mathbb{R}^p \times \mathbb{R}^p \xrightarrow{K} \mathbb{R}$ soit un noyau convenable.

Comment fabriquer des noyaux?

En préalable, deux ou trois choses à savoir sur les matrices semi-définies positives.

La première (aisée à vérifier) est qu'une combinaison linéaire à coefficients positifs de matrices semi-définies positives est semi-définie positive.

La deuxième concerne le produit terme à terme deux matrices, dit produit de Hadamard et noté $A\odot B$

Lemme du produit de Schur

Soit $A, B \in M_{n \times n}(\mathbb{R})$ semi-définies positives, alors $A \odot B$ l'est aussi.

DEMONSTRATION

Il est clair que A, B étant symétrique, $A \odot B$ l'est aussi.

A semi-définie positive est diagonalisable avec une matrice de passage orthogonale et toutes ses valeurs propres sont positives ou nulles :

$$A = Q \operatorname{Diag}(\lambda_1, \dots, \lambda_n) Q^t \quad (0 \le \lambda_i) \text{ et on } a : a_{i,j} = \sum_{k=1}^n \lambda_k q_{i,k} q_{j,k}.$$

Soit
$$U = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix} \in \mathbb{R}^n$$
, $U^t A \odot BU = \sum_{i,j=1}^n u_i a_{i,j} b_{i,j} u_j = \sum_{i,j=1}^n u_i \left(\sum_{k=1}^n \lambda_k q_{i,k} q_{j,k} \right) b_{i,j} u_j = \sum_{k=1}^n \lambda_k \sum_{i,j=1}^n u_i q_{i,k} b_{i,j} u_j q_{j,k}$

En posant :
$$U_k = \begin{pmatrix} u_1 q_{1,k} \\ \vdots \\ u_n q_{n,k} \end{pmatrix}$$
, $\sum_{i,j=1}^n u_i q_{i,k} b_{i,j} u_j q_{j,k}$ peut s'écrire : $U_k^t B U_k$. D'où :

$$U^{t}A \odot BU = \sum_{k=1}^{n} \lambda_{k} U_{k}^{t} BU_{k}$$

Or B étant semi-définie positive, $0 \leq U_k^t B U_k$ et λ_k aussi.

La troisième, conséquence immédiate du lemme précédent est que si $A = (a_{i,j})_{\substack{1 \le i \le n \\ 1 \le j \le n}}$ est semi-définie positive, la matrice $A^{\hat{q}} = (a_{i,j}^q)_{\substack{1 \le i \le n \\ 1 \le j \le n}} q \in \mathbb{N}$ l'est aussi.

On sait que $K(U,V) = \langle U;V \rangle$ est un noyau. Pour en fabriquer d'autres, on peut s'appuyer sur la proposition suivante :

Proposition

Si $K_1(U,V), K_2(U,V)$ sont des noyaux, alors:

1)
$$aK_1(U,V), +bK_2(U,V) \quad (a,b \in \mathbb{R}_+^*)$$

2)
$$aK_1(U,V)+b \quad (a,b \in \mathbb{R}_+^*)$$

3)
$$K_1(U,V) \times K_2(U,V)$$

4)
$$K_1(U,V)^k \quad (k \in \mathbb{N}_+^*)$$

5)
$$g(U)K_1(U,V)g(V)$$
 où g est une fonction de $\mathbb{R}^p \xrightarrow{g} \mathbb{R}$

6)
$$\exp(K_1(U,V))$$

sont des noyaux.

DEMONSTRATION:

Remarque: ces nouveaux noyaux et donc les matrices correspondantes sont évidemment symétriques.

On notera $G_K(X_1,\dots,X_n)$, la matrice correspondant au noyau K.

- 1) $G_{aK_1+bK_2}(X_1,\cdots,X_n) = aG_{K_1}(X_1,\cdots,X_n) + bG_{K_2}(X_1,\cdots,X_n)$. Comme $G_{K_1}(X_1,\cdots,X_n)$ et $G_{K_2}(X_1,\cdots,X_n)$ sont semi-définies positives, $G_{aK_1+bK_2}(X_1,\cdots,X_n)$ l'est aussi et $aK_1(U,V)$, $+bK_2(U,V)$ est un noyau.
- 2) C'est le cas précédent avec $\forall U, V: K_2(U, V) = 1$.

Dans ce cas
$$G_{K_2}(X_1, \dots, X_n) = \begin{pmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{pmatrix} = \vec{\mathbf{1}} \vec{\mathbf{1}}^t$$
 où $\vec{\mathbf{1}} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$. Cette matrice est semi définie positive car $\forall Z \in \mathbb{R}^n : \underline{Z}^t \vec{\mathbf{1}} \vec{\mathbf{1}}^t \underline{Z} = (Z^t \vec{\mathbf{1}})^2$.

- 3) $G_{K_1 \times K_2}(X_1, \dots, X_n) = G_{K_1}(X_1, \dots, X_n) \odot G_{K_2}(X_1, \dots, X_n)$. Il suffit d'appliquer le lemme de Schur.
- 4) Conséquence immédiate de 3).

5) Soit
$$K_{g}(U,V) = g(U)K_{1}(U,V)g(V)$$
. Alors:

$$G_{K_g}(X_1,\dots,X_n) = (g(X_i)K_1(X_i,Y_j)g(Y_j))_{\substack{1 \le i \le n \\ 1 \le j \le n}}$$

=
$$DG_{K_1}(X_1,\dots,X_n)D$$
 en notant D = $Diag(g(X_1),\dots,g(X_n))$

$$U^{t}G_{K_{o}}\left(X_{1},\cdots,X_{n}\right)U=U^{t}DG_{K_{1}}\left(X_{1},\cdots,X_{n}\right)DU$$

 $\text{Comme D diagonale $D^t=D$ et $U^tG_{K_s}\left(X_1,\cdots,X_n\right)U=\left(DU\right)^tG_{K_1}\left(X_1,\cdots,X_n\right)DU\geq 0$.}$

6) On pose : $K_1(U,V) = A = (a_{i,j})_{\substack{1 \le i \le n \\ 1 \le j \le n}}$. A est semi-définie positive.

On définit :
$$A_k = \left(\sum_{q=0}^k \frac{1}{q!} a_{i,j}^q\right)_{\substack{1 \le i \le n \\ 1 \le j \le n}} = \sum_{q=0}^k \frac{1}{q!} \left(a_{i,j}^q\right)_{\substack{1 \le i \le n \\ 1 \le j \le n}} = \sum_{q=0}^k \frac{1}{q!} A^{\hat{q}}$$

D'après les remarques préalables $A^{\hat{q}}$ est semi-définie positive et $A_k = \sum_{q=0}^k \frac{1}{q!} A^{\hat{q}}$ aussi.

 $\forall U \in \mathbb{R}^n : U^t A_k U \ge 0$. En passant à la limite $\lim_{k \to +\infty} \left(U^t A_k U \right) = U^t \lim_{k \to +\infty} \left(A_k \right) U \ge 0$.

Or
$$\lim_{k \to +\infty} A_k = \left(\lim_{k \to +\infty} \left(\sum_{q=0}^k \frac{1}{q!} a_{i,j}^q \right) \right)_{\substack{1 \le i \le n \\ 1 \le j \le n}} = \left(\exp\left(a_{i,j} \right) \right)_{\substack{1 \le i \le n \\ 1 \le j \le n}}$$

Noyaux usuels

Parmi les noyaux utilisés en pratique, on trouve :

Noyau polynomial	$K(U,V) = (coef 0 + \langle U; V \rangle)^{degree}$	Si coef0 = 0 et degree = 1, noyau linéaire
Noyau radial	$K(U,V) = \exp(-\gamma \ U - V\ ^2)$	$\gamma \in \mathbb{R}^*_{\scriptscriptstyle{+}}$ par défaut $\gamma = \frac{1}{}$
Noyau sigmoïde	$K(U,V) = \tanh(\gamma \langle U; V \rangle + coef 0)$	$p = \mathbb{R}_+$ par deraut $p = p$

2) et 4) justifient l'appellation de noyau polynomial pour : $K\left(U,V\right) = \left(coef\ 0 + \left\langle U\ ;V\right\rangle\right)^{degree}$

En ce qui concerne $K(U,V) = \exp(-\gamma ||U-V||^2)$:

$$\exp(-\gamma \|U - V\|^2) = \exp(-\gamma (U - V)^t (U - V)) = \exp(-\gamma (U^t U - 2U^t V + V^t V))$$
$$= \exp(-\gamma U^t U) \exp(2\gamma U^t V) \exp(-\gamma V^t V)$$

Or, comme $\gamma \in \mathbb{R}_+^*$, $K(U,V) = 2\gamma U'V$ est un noyau. Donc, d'après 6) $\exp(2\gamma U'V)$ aussi. Et

5) Permet de conclure avec $g(U) = \exp(-\gamma U'U)$.

Le noyau sigmoïde, bien qu'utilisé en pratique⁴ ne répond pas au critère $G\left(X_1,\cdots,X_n\right)$ semi-définie positive pour tout $\left(X_1,\cdots,X_n\right) \in \mathbb{R}^p \times \cdots \times \mathbb{R}^p$.

Soit, par exemple,
$$X_1 = \begin{pmatrix} \sqrt{2} \\ \sqrt{2} \end{pmatrix}$$
, $X_2 = \begin{pmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}$. Avec un « noyau » sigmoïde, on a :

 $G(X_1, X_2) = \begin{pmatrix} \tanh(4) & \tanh(2) \\ \tanh(2) & \tanh(1) \end{pmatrix}$ et cette matrice ayant une valeur propre (\(\times -0.0908665765\))

strictement négative n'est pas semi-définie positive.

⁴ Il est proposé parmi les options dans le package e1071 de R consacré à SVM.

Des scores aux probabilités

D'autres méthodes que SVM peuvent être utilisées devant ce genre de problème (variables explicatives quantitatives, variable à expliquer binaire) comme, par exemple la régression logistique. À l'encontre de SVM qui fournit des scores dont le signe permet la discrimination, la régression logistique fournit des probabilités d'appartenance. La méthode suivante permet la conversion des scores en probabilités.

Méthode de Platt

Elle consiste à transformer le score $f(Z) = \sum_{i \in S} \tilde{\alpha}_i y_i K(Z, X_i) + \tilde{\beta}_0$ en probabilité par une fonction sigmoïde :

$$P[y=1/Z] = \frac{1}{1 + \exp(af(Z) + b)}$$

a, b étant deux constantes à déterminer par le maximum de vraisemblance

Probabilité de l'échantillon X

Pour $i \in 1 \cdots n$, on note $p_i = P[y = 1/X_{i,\bullet}]$.

Pour la ligne i, la probabilité est donc p_i si $y_i=1,1-p_i$ si $y_i=-1$, ce qui peut se synthétiser en :

$$p_i^{\frac{y_i+1}{2}} (1-p_i)^{\frac{-y_i+1}{2}}$$
 ou encore $p_i^{[y_i=1]} (1-p_i)^{[y_i=-1]}$

La vraisemblance de l'échantillon X est donc $\prod_{i=1}^n p_i^{\frac{y_i+1}{2}} \left(1-p_i\right)^{\frac{-y_i+1}{2}}$ et en passant au logarithme, le maximum de vraisemblance est :

$$\max_{a,b} \left(\sum_{i=1}^{n} ([y_i = 1] \ln p_i + [y_i = -1] \ln (1 - p_i)) \right)$$

Avec l'exemple marge souple				a =	-0,194
Ave	c r exemple	e marge so	прте	b =	0,046
individus i	x1	x2	у	f(Xi)	P[y=1/x1 ,x2]
1	1	1	-1	-0,667	45,6%
2	4	5	-1	-1,000	44,0%
3	6	9	-1	-1,667	40,9%
4	8	7	-1	-0,333	47,2%
5	7	1	-1	1,333	55,3%
6	1	3	1	-1,333	42,5%
7	5	1	1	0,667	52,1%
8	13	6	1	1,667	56,9%
9	9	4	1	1,000	53,7%
10	12	7	1	1,000	53,7%

Ces probabilités sont peu significatives du fait d'une base d'apprentissage très étriquée avec deux grosses « erreurs » 5 et 6...

 $^{^{5}~\}mathrm{Si}~P$ est une proposition logique, $\left[P\right]\!=\!1\,\mathrm{si}\,P$ vraie $0\,\mathrm{sinon}$.

Un exemple avec des vraies données spam or not

La base de données Spambase est due à :

```
Hopkins, M., Reeber, E., Forman, G., & Suermondt, J. (1999). Spambase [Dataset]. UCI Machine Learning Repository. <a href="https://doi.org/10.24432/C53G6X">https://doi.org/10.24432/C53G6X</a> .
```

Elle est téléchargeable sur la page⁶: https://archive.ics.uci.edu/dataset/94/spambase qui contient de plus nombre d'informations à son sujet.

Cette base décrit 4601 méls à l'aide de 57 variables quantitatives en précisant le statut du mél Spam (o/n).

Parmi ces 57 variables explicatives, on trouve:

- ✓ 48 représentent la fréquence d'apparition de certains mots dans le texte du mél. Par exemple la variable wf_money indique la fréquence du mot anglais « money » (tous les méls analysés sont en anglais).
- ✓ 6 représentent la fréquence d'apparition de certains caractères spéciaux dans le texte comme #, \$...
- \checkmark Les 3 dernières sont en lien avec des séquences de plusieurs lettres écrits en majuscules. La variable Spam (o/n) est la variable à expliquer.

Après un tri aléatoire (initialement les 1873 spams figuraient en premier, les 2788 non spams en second), la base a été importée dans sous le nom de *spam* pour y être traitée à l'aide du package e1071. Puis, la base a été scindée en 2 :

- 1. Les 3601 premières lignes pour constituer la base d'apprentissage (train).
- 2. Les 1000 lignes suivantes pour constituer une base de test (test).

```
> train=spam[1:3601,]
> test=spam[3602:4601,]
> dim(train)
[1] 3601 58
> dim(test)
[1] 1000 58
```

Remarques préalables

Dans l'exposé mathématique que nous avons fait sur SVM, les vecteurs $X_{i,\bullet}$ sont étiquetés par $y_i \in \{-1,1\}$. Mais il s'agit d'une pure convention, convention justifiée car elle permet d'exprimer simplement les conditions de séparation : $\forall i \in 1 \cdots n : 1 \leq y_i \left(X_{i,\bullet}\beta + \beta_0\right)$. Mais il s'agit fondamentalement d'une discrimination binaire. La variable à expliquer Spam dans l'exemple

⁶ Elle se trouve aussi dans le classeur à télécharger : spamdata.xls

prend comme valeur o ou n. Avec le package e1071 (même chose avec Tanagra), la variable à expliquer, même si elle est codée numériquement, doit être de type factor.

Le classement s'opère sur le signe de la fonction de score f(Z) et il faut reconnaitre quelle convention a adopté le logiciel. Dans notre exemple, par anticipation, $0 \le f(Z)$ correspond à Spam = n.

La fonction sym () du package e1071 réalise le travail d'apprentissage :

Quelques explications

Arguments de la fonction

- ✓ formula=Spam~. précise la variable à expliquer.
- ✓ data=train indique la data-frame à utiliser.
- \checkmark cost = 2 fixe la valeur de la constante C (Γ dans notre texte) qui règle la tolérance aux erreurs.
- ✓ kernel = "radial" indique le noyau à utiliser à choisir parmi :
 "linear", "polynomial", "radial", "sigmoid".
- ✓ probability = TRUE calcul des coefficients a et b intervenant dans la formule de conversion score → probabilité.
- ✓ scale=TRUE pour standardiser les variables.

Résultats

Les résultats sont englobés dans l'objet m1. La commande print (m1) n'en dévoile qu'une petite partie. La commande attributes (m1) révèle d'autres attributes :

```
names
"call"
             "type"
                            "kernel"
                                           "cost"
                                                         "degree"
                                                                              "gamma"
"coef0"
             "nu"
                            "epsilon"
                                           "sparse"
                                                         "scaled"
                                                                              "x.scale"
"y.scale"
                                           "tot.nSV"
                                                         "nSV"
             "nclasses"
                            levels"
                                                                              "labels"
"SV"
             "index"
                            "rho"
                                                         "probA"
                                                                              "probB"
                                           "compprob"
"sigma"
                                                         "decision.values"
             "coefs"
                            "na.action"
                                          "fitted"
                                                                              "terms"
```

Parmi ces attributs:

- ✓ m1\$SV fournit la liste des vecteurs supports avec leurs coordonnées.
- \checkmark m1\$coef donne pour chaque vecteur support $\tilde{\alpha}_i$ y_i .

Avec un noyau "linear", $\tilde{\beta}$ peut être obtenue par : beta =t (m1\$coefs) %*% m1\$SV.

- \checkmark m1\$rho donne la valeur de $\tilde{\beta}_0$.
- ✓ m1\$decision.values donne la liste pour $i \in 1 \cdots n$ des $f(X_{i,\bullet})$.
- ✓ m1\$probA, m1\$probB les valeurs des coefficients a et b évoqués en amont.

Prédiction

Après l'apprentissage vient la phase de prédiction. Elle se fait par la fonction predict (). Dans sa version courte, elle possède deux arguments :

- 1. Un objet retourné par la fonction SVM()
- 2. Une data-frame de même structure que celle qui a servi à l'apprentissage.

Avant de la mettre en œuvre sur la base test, on peut commencer par le faire sur la base train :

```
> p1=predict(m1, train)
> p1[1:10]
1 2 3 4 5 6 7 8 9 10
n n o o o n n n o o
Levels: n o
```

On peut facilement créer la matrice de confusion qui confronte valeur réelle et valeur prédite :

Pour la base test, on obtient:

Résumé

Matrice de confusion		Prédiction		
train		0	n	
Valeur réelle	0	2119	59	
valeur reelle	n	101	1322	
	Σ	2220	1381	

_		
2178	TVP (sensibilité)	97,3%
1423	TVN (spécificité)	92,9%
3601	Taux d'erreur	4,4%

Matrice de confusion		Prédiction		
test		0	n	
Valeur réelle	0	585	25	
valeur reelle	n	50	340	
	Σ	635	365	

Σ		
	TVP (sensibilité)	95,9%
390	TVN (spécificité)	87,2%
1000	Taux d'erreur	7,5%

On observe un taux d'erreurs plus élevé sur la base test. C'est un constat courant et assez logique, puisque avec la base train ce sont les mêmes données qui servent à l'apprentissage et que la fonction score f est élaboré pour « coller » au mieux à ces données.

La sensibilité dans les deux cas est assez bonne : $\sim 3-4\%$ de spams non détectés.

En revanche la spécificité est mauvaise sur la base test $\sim 13\%$ de méls sains sont classés en spam.

La fonction predict () dans sa version longue s'écrit en se limitant pour des raisons évidentes aux 5 derniers méls :

```
> predict(m1, test[996:1000,],probability=TRUE, decision.values=TRUE)
4597 4598 4599 4600 4601
             n
       n
attr(, "decision.values")
            n/o
     0.9998383
4597
4598
     0.7109559
4599
     0.1634921
4600 -1.5943902
4601 -1.3818685
attr(,"probabilities")
              n
4597 0.92451892 0.07548108
4598 0.85506814 0.14493186
4599 0.59642450 0.40357550
4600 0.01704965 0.98295035
4601 0.02883121 0.97116879
Levels: n o
```

On observe que les valeurs positives de l'attribut decision.values correspondent à des méls « sains » (Spam=n).

La pratique de SVM est un art subtil. Deux choix cruciaux sont à effectuer :

- 1. Le choix du noyau.
- 2. Le choix du paramètre C (cost).

Le package e1071 propose la fonction tune () pour guider l'utilisateur dans son choix :

```
>tune(svm, Spam~., data= train,
  ranges =list(kernel = c('linear', 'polynomial', 'radial'),
  cost =c(0.5,1,2,5,10)), tunecontrol = tune.control(sampling ="cross"))

Parameter tuning of 'svm':
  - sampling method: 10-fold cross validation
  - best parameters:
  kernel cost
  radial 5
  - best performance: 0.06414589
```

La méthode employée se fonde sur une validation croisée. Les données sont scindées en 10 morceaux (folds) et tour à tour chacun des morceaux sert de test après apprentissage sur les 9 autres. Et il faut itérer la procédure pour les 4 costs et les 3 kernels proposés. De ce fait le temps d'exécution est assez long... (plus de 4 minutes sur mon PC).

Nous avions fait le bon choix pour le noyau (radial), mais pas pour cost. Si on recommence les calculs avec cost=5, on obtient les matrices de confusion suivantes :

Matrice de confusion		Prédiction	
train		0	n
Valeur réelle	0	2135	43
valeur reelle	n	89	1334
	Σ	2224	1377

2		
2178	TVP (sensibilité)	98,0%
1423	TVN (spécificité)	93,7%
3601	Taux d'erreur	3,7%

Matrice de confusion		Prédiction	
test		0	n
Valeur réelle	0	585	25
valeur reelle	n	46	344
	Σ	631	369

_			
	610	TVP (sensibilité)	95,9%
	390	TVN (spécificité)	88,2%
-	1000	Taux d'erreur	7,1%

où on peut constater une diminution des taux d'erreur.

Notations

Notations générales

E	Nombre d'éléments	d'un	ensemble	fini	Е.

$$\mathcal{P}(E)$$
 L'ensemble des parties de l'ensemble E .

$$k \cdots n$$
 Ensemble des entiers $k \le i \le n$

$$[P] P ext{ étant une proposition } [P] = \begin{cases} 1 ext{ si P est vraie} \\ 0 ext{ sinon} \end{cases}$$

Algèbre linéaire

$$\mathcal{M}_{n \times p}(\mathbb{R})$$
 L'espace des matrices réelles de n lignes et p colonnes

$$\mathcal{M}_n(\mathbb{R})$$
 L'espace des matrices réelles carrées de n lignes et n colonnes

$$(a_{i,j})_{\substack{i=1\cdots n\\j=1\cdots p}}$$
 Matrice de taille $n\times p$ ayant pour terme $a_{i,j}$ en ligne i , colonne j .

$$a_{i,j}$$
 $A_{i,j}$ Le terme de la matrice A situé en ligne i , colonne j

$$A_{i,\bullet}$$
 La ligne i de la matrice A .

$$A_{\cdot,j}$$
 La colonne j de la matrice A .

$$A^{t}$$
 Matrice transposée de la matrice A .

$$\vec{\mathbf{1}}_{n}, \vec{\mathbf{1}}$$
 Le vecteur de \mathbb{R}^{n} dont toutes les composantes sont égales à 1.

$$E_i$$
 Le $i^{\grave{e}me}$ vecteur de la base canonique de \mathbb{R}^n .

$$E^*$$
 Le dual de l'espace vectoriel E .

$$\left\langle U\,;V\right\rangle$$
 Produit scalaire de 2 vecteurs d'un espace euclidien.

$$||U||$$
 Norme euclidienne.

$$A \leq B \; ; \; A,B \in \mathcal{M}_{n \times p}(\mathbb{R}) \qquad \forall i \in 1 \cdots n, \, \forall j \in 1 \cdots p \; : \; a_{i,j} \leq b_{i,j}$$

$$A^{\hat{q}} = \left(a_{i,j}^{q}\right)_{\substack{1 \le i \le n \\ 1 \le j \le n}} q \in \mathbb{N}$$

$$Analyse:pour \mathbb{R}^{p} \xrightarrow{h} \mathbb{R}$$

$$\nabla(h) \qquad \text{Gradient} : \nabla(h) = \begin{pmatrix} \frac{\partial h}{\partial x_1} \\ \vdots \\ \frac{\partial h}{\partial x_n} \end{pmatrix}$$

$$\nabla^{2}(h) \qquad \text{Matrice symétrique hessienne } \left(\frac{\partial^{2} h}{\partial x_{i} \partial x_{j}}\right)_{\substack{i=1\cdots p\\j=1\cdots p}}$$

Notations pour SVM

(On s'est efforcé d'adopter les notations courantes dans ce domaine)

 $p \in \mathbb{N}^*$ Nombre de variables explicatives numériques.

 $n \in \mathbb{N}^*$ Nombre de cas.

 $X \in M_{n \times p}(\mathbb{R})$ Matrice des données.

 $y \in \{-1, 1\}^n$ La variable binaire à expliquer.

 $\Gamma \in \mathbb{R}_{+}$ La constante cost

 $\tilde{\beta}_0 \in \mathbb{R}, \, \tilde{\beta} \in \mathbb{R}^p, \, \tilde{\xi} \in \mathbb{R}^n_+$ Solution du primal. $\tilde{\xi}$ variables d'écart (slack variables)

 $\tilde{\alpha} \in \mathbb{R}^n_+$ Solution du dual.

S vecteurs supports $S = \{i \in \mathbb{N}^* / 0 < \tilde{\alpha}_i \}$

 $S_{I} \text{ type } I$ $S_{I} = \{i \in \mathbb{N}^{*} / 0 < \tilde{\alpha}_{i} < \Gamma\}$

 S_{II} type II $S_{II} = \left\{ i \in \mathbb{N}^* / \tilde{\alpha}_i = \Gamma \right\}$

La fonction score:

$$\left(\mathbb{R}^{p}\right)^{*} \xrightarrow{f} \mathbb{R}$$

$$f(Z) = \begin{cases} Z\tilde{\beta} + \beta_{0} \text{ si noyau linéaire} \\ \sum_{i \in S} \tilde{\alpha}_{i} y_{i} K\left(Z, X_{i, \bullet}\right) + \tilde{\beta}_{0} \text{ sinon} \end{cases}$$

Classeurs Excel à télécharger

- ✓ ex1 dual.xlsx
- √ ex2 effet cost.xlsx
- ✓ spamdata.xls

Quelques références

- 1 Chang C.C., Lin C.J., LIBSVM: a library for support vector machines https://www.csie.ntu.edu.tw/%7Ecjlin/libsvm/
- 2 Frédéric Meunier Introduction à l'optimisation https://cermics.enpc.fr/~meuniefr/IntroOptim.pdf
- 3 Freie Universität Berlin Support Vector Machines Algorithm

 https://www.geo.fu-berlin.de/en/v/soga-r/Machine-Learning/SVM/index.html
- 4 Jacques Cellier Algèbre Linéaires des bases aux applications Presses Universitaires de Rennes 2008
- 5 Martin Haugh Support Vector Machines (and the Kernel Trick)
 http://www.columbia.edu/~mh2078/MachineLearningORFE/SVMs MasterSlides.pdf
- 6 P.G. Ciarlet Introduction à l'analyse numérique matricielle et à l'optimisation MASSON 1994
- 7 Ricco Rakotomalala Normalisation des scores
 https://eric.univ-lyon2.fr/ricco/cours/slides/calibration.pdf
- 8 Ricco Rakotomalala SVM Support Vector Machine
 https://eric.univ-lyon2.fr/ricco/cours/slides/svm.pdf
- 9 Shai Shalev-Shwartz *IFT-7002 Fondements de l'apprentissage machine SVM et méthodes à noyaux Traduit et adapté par Mario Marchand Université Laval*http://www2.ift.ulaval.ca/~mmarchand/IFT7002/SVMetNoyaux.pdf
- Vladimir N. Vapnik The Nature of Statistical Learning Theory

 https://statisticalsupportandresearch.wordpress.com/wp-content/uploads/2017/05/vladimir-vapnik-the-nature-of-statistical-learning-springer-2010.pdf

Vladimir N. Vapnik est un des inventeurs de SVM

Table des matières

SVM, bases mathématiques	
Optimisation quadratique	
Karush-Kuhn-Tucker	
Lagrangien	
Retour au problème quadratique	
Application à SVM	5
Deux nuages linéairement séparables, marge rigide	
Deux nuages non linéairement séparables, marge souple	9
Séparation non linéaire	16
L'astuce du noyau (kernel trick) en détail	
Quelles sont les propriétés requises pour un noyau ?	17
Comment fabriquer des noyaux ?	19
Noyaux usuels	
Des scores aux probabilités	22
Un exemple avec des vraies données spam or not	23
Quelques explications	
Prédiction	
Notations	28
Quelques références	30