

# Apprentissage Machine / Statistique

## *Support Vector Machines (SVM)*

PHILIPPE BESSE

INSA de Toulouse  
Institut de Mathématiques

## Principes généraux

- Séparateur à Vaste Marge (SVM)
- Machine à Vecteurs Support (MVS)
- Apprentissage en discrimination :  $\{-1, 1\}$
- Etendu à  $m > 2$  et  $\mathbb{R}$
- **Hyperplan** de marge optimale pour la généralisation
- Vapnik (1998) et **VC**-dimension
- Contrôle de la **complexité**
- L'objectif, seulement l'objectif
- **Coût calcul** fonction de  $n$ , pas de  $p$

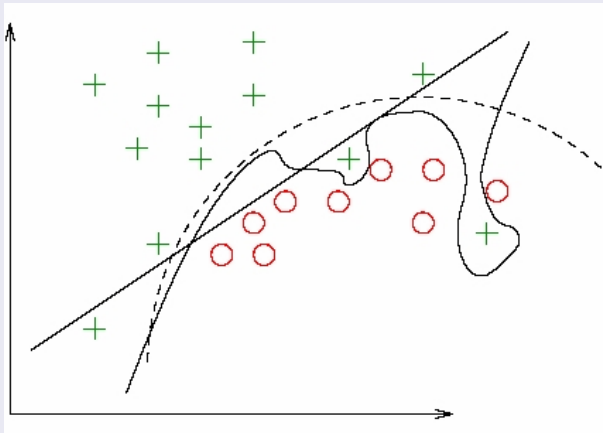
## Spécificités

- Ramener la discrimination à un problème **linéaire**
- Problème d'optimisation **sous-contrainte** et support
- Utilisation d'un espace **intermédiaire** (*feature space*)
- Produit scalaire et **noyau** reproduisant

## Remarques

- Efficacité et **flexibilité** des noyaux
- Schölkopf et Smola (2002)
- `www.kernel-machines.org`

## Sur-ajustement



*Frontière, complexité, généralisation et VC-dimension*

## Notations

- $Y$  à valeurs dans  $\{-1, 1\}$
- $X = X^1, \dots, X^p$  les variables prédictives
- $Y = f(X)$  un modèle pour  $Y$
- Un échantillon statistique de loi  $F$

$$\mathbf{z} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$$

- Estimation de  $\hat{f}$  de  $f$ ,  $(\mathbb{R}^p \text{ (ou } \mathcal{F}) \mapsto \{-\infty, \infty\})$
- par minimisation de :

$$P(f(X) \neq Y)$$

## Définition de la marge

- $f$  définie par une fonction réelle  $f : \hat{f} = \text{signe}(f)$
- L'erreur devient :  $P(f(\mathbf{X}) \neq Y) = P(Yf(\mathbf{X}) \leq 0)$
- $|Yf(\mathbf{X})|$  est un indicateur de confiance
- $Yf(\mathbf{X})$  est la **marge** de  $f$  en  $(\mathbf{X}, Y)$

## Espace hilbertien

- $\Phi : \mathbb{R}^p (\text{ou } \mathcal{F}) \mapsto \mathcal{H}$
- $\mathcal{H} : \text{feature space}$  de grande dimension avec produit scalaire
- $\Phi$  ramène à un problème linéaire : hyperplan séparateur
- Première approche :  $\Phi$  est la fonction identité

## Recherche du plan de marge maximale

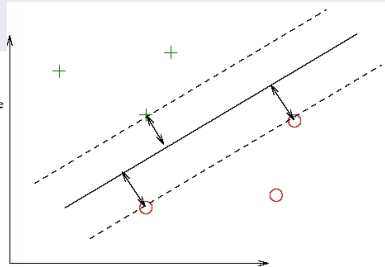
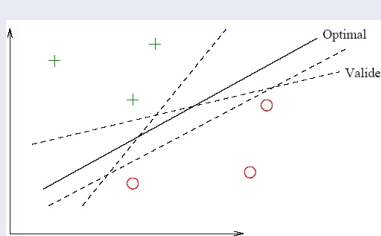
- Un hyperplan est défini à l'aide du produit scalaire de  $\mathcal{H}$  :

$$\langle \mathbf{w}, \mathbf{x} \rangle + b = 0$$

où  $\mathbf{w}$  est un vecteur orthogonal au plan

- Le signe de la fonction  $f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$  indique la position de  $\mathbf{x}$  à prédire
- Un point est bien classé si et seulement si :  $yf(\mathbf{x}) > 0$
- $(\mathbf{w}, b)$  est défini à un coef. près ; on impose :  $yf(\mathbf{x}) \geq 1$
- Un plan  $(\mathbf{w}, b)$  est un séparateur si :  $\forall i \quad y_i f(\mathbf{x}_i) \geq 1$
- Distance de  $\mathbf{x}$  au plan  $(\mathbf{w}, b)$  :  $d(\mathbf{x}) = \frac{|\langle \mathbf{w}, \mathbf{x} \rangle + b|}{\|\mathbf{w}\|} = \frac{|f(\mathbf{x})|}{\|\mathbf{w}\|}$
- La marge du plan a pour valeur :  $\frac{2}{\|\mathbf{w}\|^2}$

## Plan de marge maximale





## Problème primal d'optimisation sous contraintes

$$\begin{cases} \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{avec } \forall i, y_i < \mathbf{w}, \mathbf{x}_i > +b \geq 1 \end{cases}$$

## Problème dual avec multiplicateurs de Lagrange

- La solution est un **point-selle**  $(\mathbf{w}^*, b^*, \boldsymbol{\lambda}^*)$  du lagrangien :

$$L(\mathbf{w}, b, \boldsymbol{\lambda}) = 1/2 \|\mathbf{w}\|_2^2 - \sum_{i=1}^n \lambda_i [y_i (< \mathbf{w}, \mathbf{x}_i > +b) - 1]$$

- Ce point-selle vérifie :  $\forall i \quad \lambda_i^* [y_i (< \mathbf{w}^*, \mathbf{x}_i > +b^*) - 1] = 0$
- Vecteurs support* :  $\mathbf{x}_i$  avec contrainte active
- Appartiennent au plan :  $y_i (< \mathbf{w}^*, \mathbf{x}_i > +b^*) = 1$

## Formule duale du lagrangien

- Plan optimal :  $\mathbf{w}^* = \sum_{i=1}^n \lambda_i^* y_i \mathbf{x}_i$  et  $\sum_{i=1}^n \lambda_i^* y_i = 0$
- $W(\boldsymbol{\lambda}) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Le point-selle maximise  $W(\boldsymbol{\lambda})$  avec  $\lambda_i \geq 0 \quad \forall i$
- Problème d'optimisation quadratique de taille  $n$
- Hyperplan optimal :  $\sum_{i=1}^n \lambda_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^* = 0$
- avec  $b^* = -\frac{1}{2} [\langle \mathbf{w}^*, \mathbf{sv}_{class+1} \rangle + \langle \mathbf{w}^*, \mathbf{sv}_{class-1} \rangle]$
- La prévision de  $\mathbf{x}$  est fournie par le signe de

$$f(\mathbf{x}) = \sum_{i=1}^n \lambda_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^*$$

## Cas non séparable

- Assouplissement des contraintes
- les termes d'erreur  $\xi_i$  contrôlent le dépassement :

$$y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq +1 - \xi_i \quad \forall i \in \{1, \dots, n\}$$

- La prédiction de  $\mathbf{x}_i$  est fautive à un vecteur si  $\xi_i > 1$
- La somme des  $\xi_i$  est une borne du nombre d'erreurs
- Nouveau problème de minimisation avec pénalisation par le dépassement de la contrainte :

$$\begin{cases} \min \frac{1}{2} \|\mathbf{w}\|^2 + \delta \sum_{i=1}^n \xi_i \\ \forall i, y_i \langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq +1 - \xi_i \end{cases}$$

## Remarques

- $\delta$  contrôle le compromis entre ajustement et généralisation
- Même forme duale mais avec les  $\lambda_i$  bornés par  $\delta$
- $n$  grand : algorithmes avec **décomposition** de l'ensemble d'apprentissage
- Capacité de généralisation dépend du nombre de **vecteurs supports** mais pas de la taille de l'espace
- Si les  $X$  sont dans une boule de rayon  $R$ , l'ensemble des hyperplans de marge fixée  $\delta$  a une VC-dimension bornée par  $\frac{R^2}{\delta^2}$  avec  $\|w\| \leq R$
- Bornes d'erreur estimables mais trop pessimistes

## Produit scalaire et noyau

- $\Phi : \mathbb{R}^p (\text{ou } \mathcal{F}) \mapsto \mathcal{H}$
- $\mathcal{H}$  muni d'un produit scalaire et de plus grande dimension
- Le problème de minimisation et la solution :

$$f(\mathbf{x}) = \sum_{i=1}^n \lambda_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^*$$

font intervenir  $\mathbf{x}$  et  $\mathbf{x}'$  par l'intermédiaire de **produits scalaires** :

$$\langle \mathbf{x}, \mathbf{x}' \rangle$$

## Astuce du noyau

- Il est inutile d'explicitier  $\Phi$
- Il suffit de calculer les produits scalaires dans  $\mathcal{H}$
- Fonction noyau  $k : \mathbb{R}^p \times \mathbb{R}^p \mapsto \mathbb{R}$  symétrique :

$$k(\mathbf{x}, \mathbf{x}') = \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle_{\mathcal{H}}$$

- Le **noyau** définit une notion de **distance**

## Exemple trivial

- $\mathbf{x} = (x_1, x_2)$  dans  $\mathbb{R}^2$
- $\Phi(\mathbf{x}) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$
- $\mathcal{H}$  de dimension 3 et de produit scalaire :

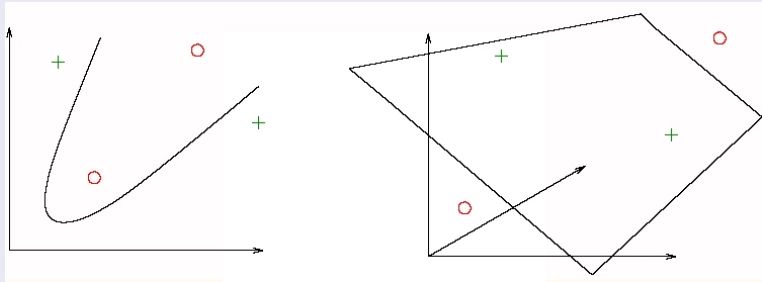
$$\begin{aligned}\langle \Phi(\mathbf{x}), \Phi(\mathbf{x}') \rangle &= x_1^2 x_1'^2 + 2x_1 x_2 x_1' x_2' + x_2^2 x_2'^2 \\ &= (x_1 x_1' + x_2 x_2')^2 \\ &= \langle \mathbf{x}, \mathbf{x}' \rangle^2 \\ &= k(\mathbf{x}, \mathbf{x}')\end{aligned}$$

## En général

- Le **produit scalaire** dans  $\mathcal{H}$  ne nécessite pas d'explicitier  $\Phi$
- Le **plongement dans  $\mathcal{H}$**  peut rendre possible la séparation linéaire



## Feature space



*Rôle de l'espace intermédiaire dans la séparation des données*

## Définition

Une fonction  $k(.,.)$  symétrique est un noyau si, pour tous les  $x_i$  possibles, la matrice de terme général  $k(x_i, x_j)$  est une matrice définie positive

- Elle définit une matrice de produit scalaire
- Dans ce cas, 'il existe un espace  $\mathcal{H}$  (Hilbert à noyau reproduisant) et une fonction  $\Phi$  tels que :

$$k(x, x') = \langle \Phi(x), \Phi(x') \rangle$$

## Attention

Condition d'existence, pas constructive et difficile à vérifier

## Noyaux classiques

- Linéaire

$$k(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$$

- Polynômial

$$k(\mathbf{x}, \mathbf{x}') = (c + \langle \mathbf{x}, \mathbf{x}' \rangle)^d$$

- Radial gaussien

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}}$$

## Noyaux spécifiques

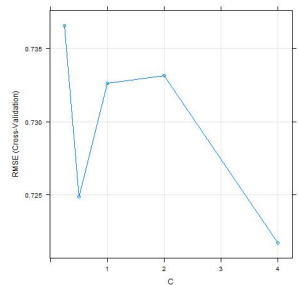
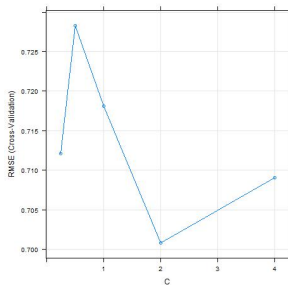
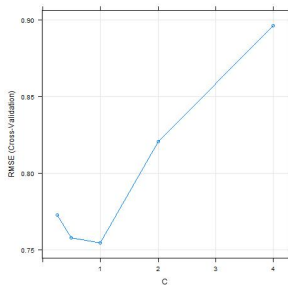
- Travail : construction d'un noyau **adapté** : reconnaissance de séquences, de caractères, l'analyse de textes, de graphes...
- Grande **flexibilité** entraîne une bonne efficacité
- **Choix** de noyau, des paramètres par validation croisée
- Paradoxe : les SVM à noyaux gaussiens dans le cas séparable ou à pénalité variable, dont de **VC-dimension** infinie

## Cas de la régression

- $Y$  est quantitative
- La fonction se décompose :  $f(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^{\infty} w_i v_i(\mathbf{x})$
- Fonction coût issue de la robustesse :

$$E(\mathbf{w}, \gamma) = \frac{1}{n} \sum_{i=1}^n |y_i - f(\mathbf{x}_i, \mathbf{w})|_{\epsilon} + \gamma \|\mathbf{w}\|^2$$

- $|\cdot|_{\epsilon}$  fonction paire, continue, identiquement nulle sur  $[0, \epsilon]$  et qui croît linéairement sur  $[\epsilon, +\infty]$
- $\gamma$  contrôle l'ajustement
- Même principe de résolution
- Noyaux de splines ou encore noyau de Dérictet

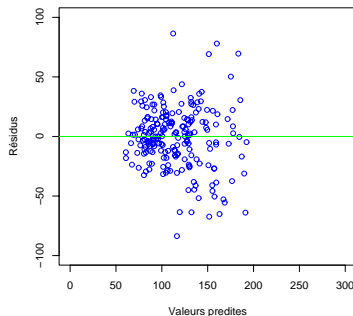
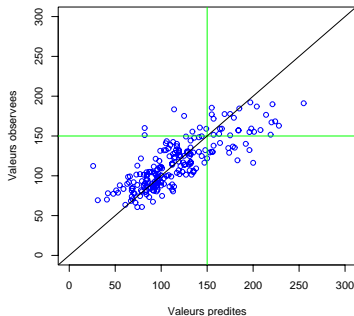


*Cookies : optimisation des SVM avec noyau linéaire*

## Exemple de discrimination

	Cancer du sein	
	benign	malignant
benign	83	1
malignant	3	50
Taux de 3%		

	Dépassement du seuil d'ozone	
	FALSE	TRUE
FALSE	161	13
TRUE	7	27
Taux de 9,6% (régression) et 12% (discrimination)		



*Ozone : Valeurs observées et résidus du test en fonction des valeurs prédites*