

# Support Vector Machines

Machine Learning

Cédric RICHARD

Université Nice Sophia Antipolis

# PRINCIPES MRE ET MRS

---

**Principe MRE.** Il s'agit de minimiser la fonctionnelle de risque

$$P_e(d) = \int \frac{1}{2} |y - d(\mathbf{x}; \mathbf{w}, b)| p(\mathbf{x}, y) d\mathbf{x} dy.$$

La densité  $p(\mathbf{x}, y)$  étant inconnue, la minimisation de  $P_e(d)$  se traduit par celle du risque empirique

$$P_{emp}(d) = \frac{1}{2n} \sum_{i=1}^n |y_i - d(\mathbf{x}_i; \mathbf{w}, b)|$$

calculable sur les données constituant l'ensemble d'apprentissage  $\mathcal{A}_n$ .

**Principe MRS.** Avec une probabilité  $1 - \varepsilon$ , on a :

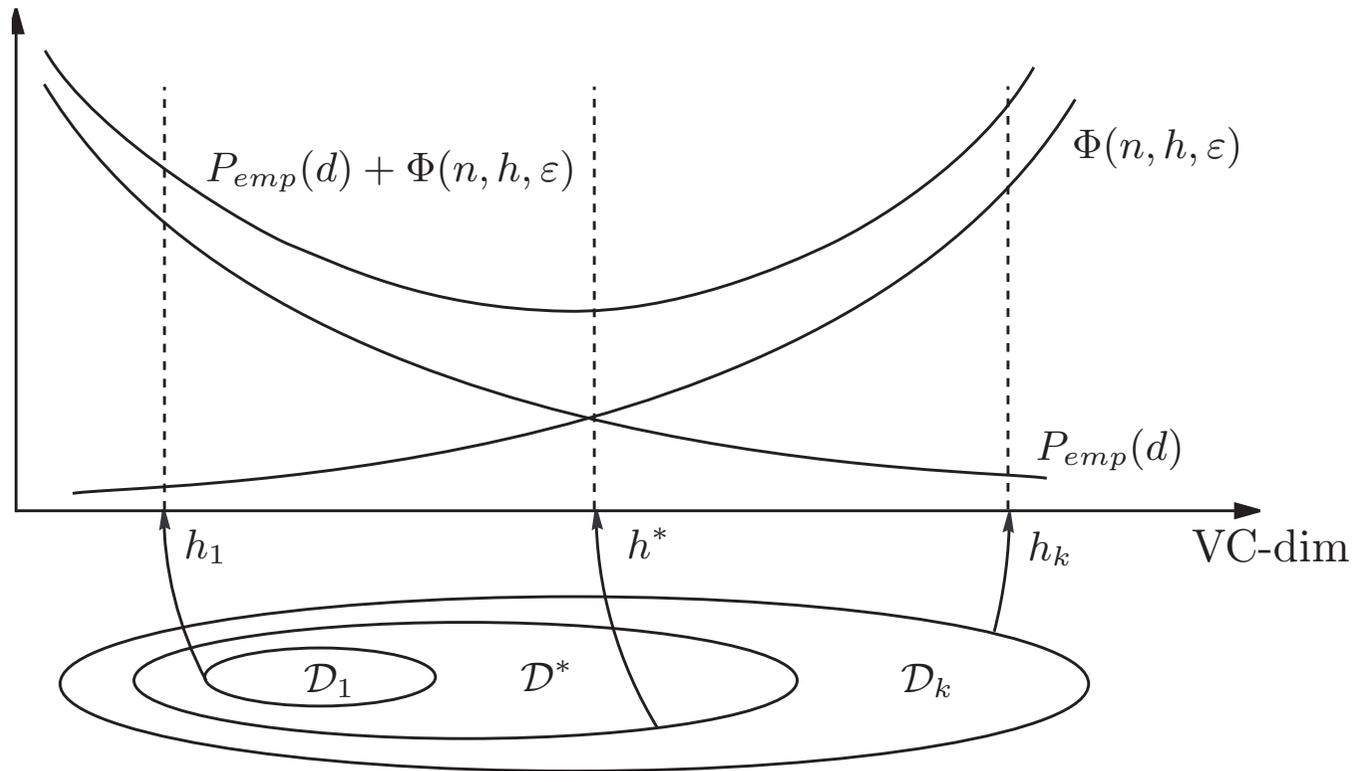
$$P_e(d) \leq P_{emp}(d) + \sqrt{\frac{h \ln \left( \frac{2n}{h} + 1 \right) - \ln \frac{\varepsilon}{4}}{n}}.$$

Le principe MRS suggère de minimiser par rapport à  $h$  le risque garanti défini par la borne supérieure  $P_{emp}(d) + \Phi(n, h, \varepsilon)$ , plutôt que le risque empirique.

# LE PRINCIPE MRS

Illustration

---



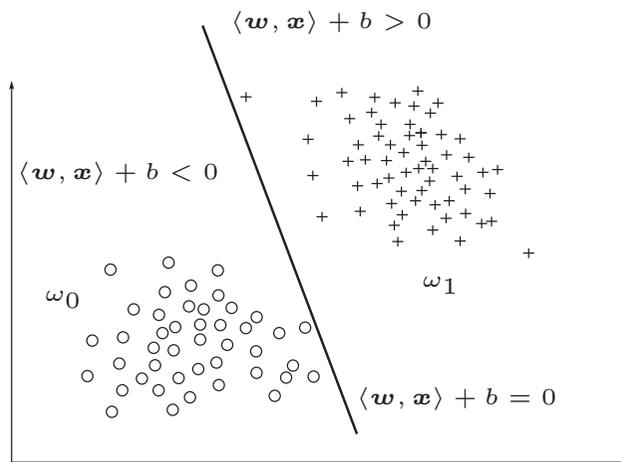
# L'ALGORITHME DU PERCEPTRON

---

L'algorithme du Perceptron vise à produire une solution d'erreur d'apprentissage minimum en minimisant le risque empirique suivant :

$$(\mathbf{w}^*, b^*) = \arg \min_{(\mathbf{w}, b)} \sum_{i=1}^n |y_i - d(\mathbf{x}_i; \mathbf{w}, b)|.$$

- ▷ Pourquoi la solution obtenue aurait-elle les meilleures performances ?
- ▷ Est-ce que la minimisation de l'erreur empirique est une bonne idée ?
- ▷ Existe-t-il une alternative ?



# CLASSIFIEURS LINÉAIRES

## Définition

---

On considère un problème de classification à 2 classes de  $n$  données dans  $\mathbb{R}^l$ , étant donné un ensemble d'apprentissage

$$\mathcal{A}_n = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}.$$

On pose  $y_i = (-1)$  si  $\mathbf{x}_i \in \omega_0$ , et  $y_i = (+1)$  si  $\mathbf{x}_i \in \omega_1$ .

Un **classifieur linéaire** est défini par

$$d(\mathbf{x}; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b).$$

# SÉPARABILITÉ LINÉAIRE

## Définition

---

On considère un problème de classification à 2 classes de  $n$  données dans  $\mathbb{R}^l$ , étant donné un ensemble d'apprentissage

$$\mathcal{A}_n = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}.$$

On pose  $y_i = (-1)$  si  $\mathbf{x}_i \in \omega_0$ , et  $y_i = (+1)$  si  $\mathbf{x}_i \in \omega_1$ .

L'équation d'un hyperplan est définie à une constante multiplicative près :

$$\langle \mathbf{w}, \mathbf{x} \rangle + b = 0 \quad \Longleftrightarrow \quad \langle \gamma \mathbf{w}, \mathbf{x} \rangle + \gamma b = 0, \quad \gamma \in \mathbb{R}^*$$

Les classes  $\omega_0$  et  $\omega_1$  sont dites **linéairement séparables** s'il existe  $\mathbf{w}$  et  $b$  tels que

$$\langle \mathbf{w}, \mathbf{x}_i \rangle + b \geq +1 \quad \forall \mathbf{x}_i \in \omega_1$$

$$\langle \mathbf{w}, \mathbf{x}_i \rangle + b \leq -1 \quad \forall \mathbf{x}_i \in \omega_0$$

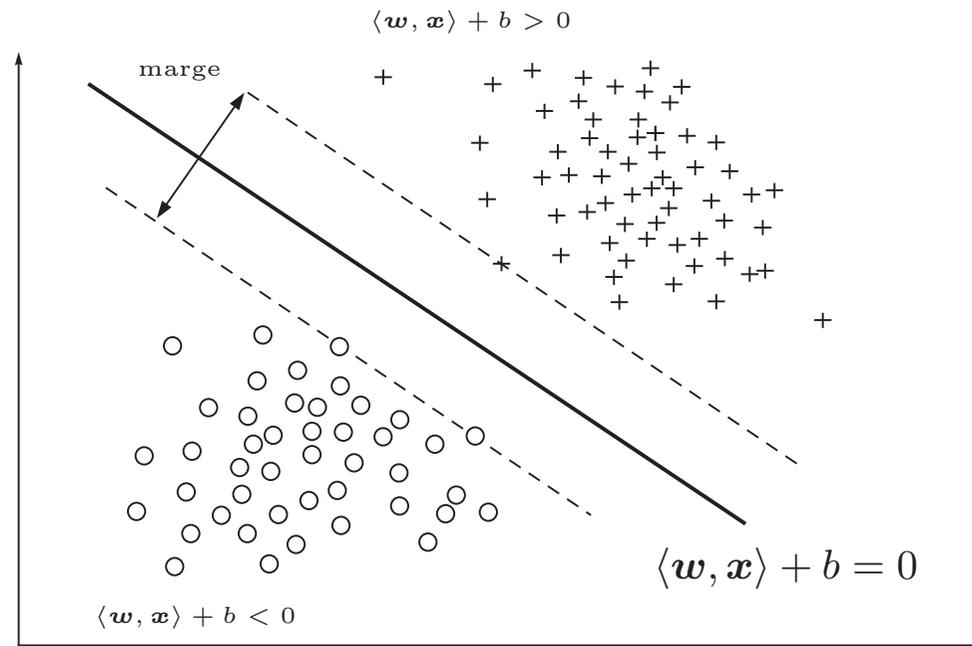
Dans la suite, on résume ce critère de séparabilité ainsi

$$y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 \geq 0 \quad \forall (\mathbf{x}_i, y_i) \in \mathcal{A}_n$$

# UN NOUVEAU PRINCIPE D'INDUCTION

---

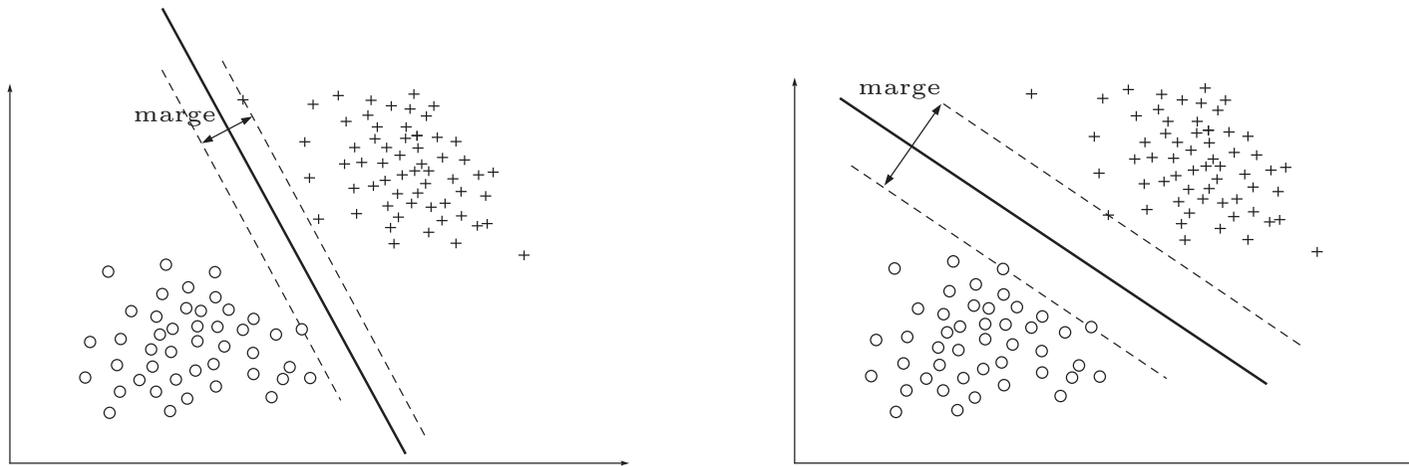
Parmi les séparateurs conduisant à une erreur empirique minimum, il convient de choisir celui de **marge maximum** (Vapnik 1965, 1992).



# UN NOUVEAU PRINCIPE D'INDUCTION

## Illustration

---



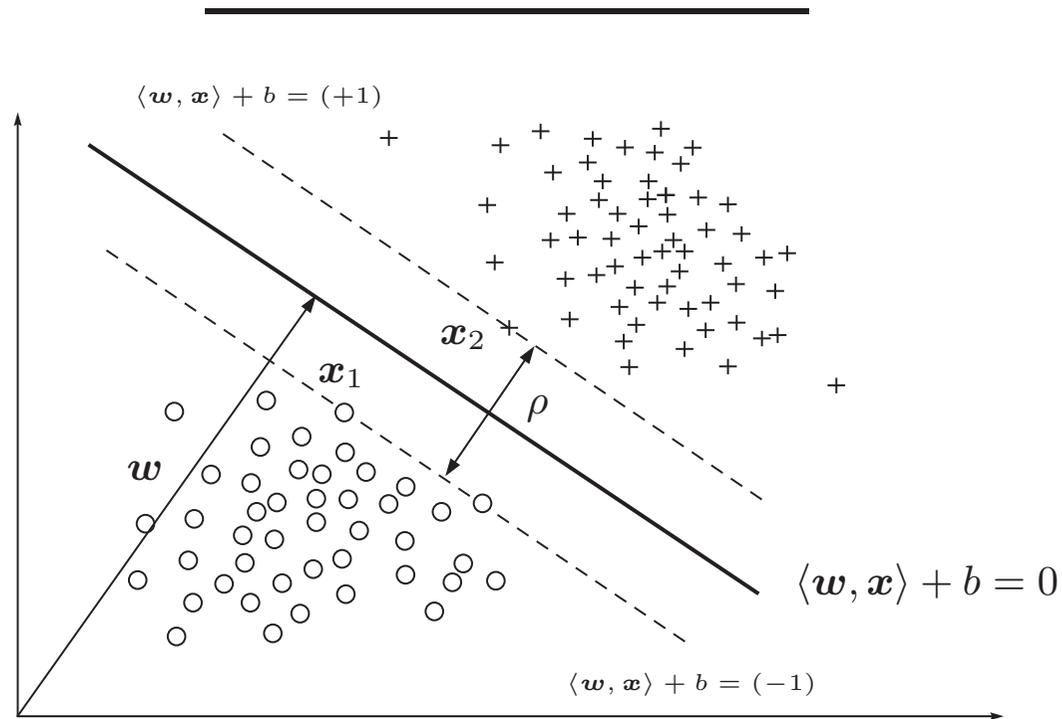
**Faible marge** : probablement de faibles performances en généralisation

**Large marge** : probablement de bonnes performances en généralisation

Ceci va être justifié de manière plus rigoureuse à présent.

# UN NOUVEAU PRINCIPE D'INDUCTION

Calcul de la marge



On a  $\langle \mathbf{w}, \mathbf{x}_2 \rangle + b = (+1)$  et  $\langle \mathbf{w}, \mathbf{x}_1 \rangle + b = (-1)$ . Il en résulte directement que

$$\rho = \left\langle \frac{\mathbf{w}}{\|\mathbf{w}\|}, \mathbf{x}_2 - \mathbf{x}_1 \right\rangle = \frac{2}{\|\mathbf{w}\|}$$

# UN NOUVEAU PRINCIPE D'INDUCTION

## Maximisation de la marge

---

On justifie la maximisation de la marge  $\rho$ , principe de base des SVM, par le résultat suivant issu de la théorie statistique de l'apprentissage.

**Théorème 1.** *Considérons les hyperplans de la forme  $\langle \mathbf{w}, \mathbf{x} \rangle = 0$ , où  $\mathbf{w}$  est normalisé de sorte qu'ils soient sous forme canonique par rapport à  $\mathcal{A}_n$ , soit :*

$$\min_{\mathbf{x} \in \mathcal{A}_n} |\langle \mathbf{w}, \mathbf{x} \rangle| = 1.$$

*L'ensemble des fonctions de décision  $d(\mathbf{x}; \mathbf{w}) = \text{sgn}\langle \mathbf{w}, \mathbf{x} \rangle$  définies à partir de  $\mathcal{A}_n$  et vérifiant la contrainte  $\|\mathbf{w}\| \leq \Lambda$  à une VC-dimension  $h$  vérifiant :*

$$h \leq R^2 \Lambda^2,$$

*où  $R$  est le rayon de la plus petite sphère centrée sur l'origine contenant  $\mathcal{A}_n$ .*

En conséquence, plus  $\rho = 2/\|\mathbf{w}\|$  est grand, plus  $h$  est petit.

# SUPPORT VECTOR MACHINES (HARD MARGIN)

## Formulation du problème d'optimisation

---

Maximiser la marge, définie par  $\rho = \frac{2}{\|\mathbf{w}\|}$ , est équivalent à minimiser  $\|\mathbf{w}\|^2$ . On implémente le principe MRS en résolvant le problème d'optimisation suivant :

Minimiser  $\frac{1}{2}\|\mathbf{w}\|^2$

sous les contraintes  $y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1, \quad 1 \leq i \leq n.$

**Remarque.** Cette formulation est valide pour des classes linéairement séparables.

# SUPPORT VECTOR MACHINES (HARD MARGIN)

## Rappels sur les multiplicateurs de Lagrange

---

La minimisation d'une fonction convexe  $f(\mathbf{x})$  sous les contraintes  $g_i(\mathbf{x}) \leq 0$ ,  $i = 1, \dots, n$ , est équivalente à la recherche du point selle du Lagrangien

$$L(\mathbf{x}; \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_{i=1}^n \alpha_i g_i(\mathbf{x}).$$

Le minimum est recherché par rapport à  $\mathbf{x}$ . Le maximum l'est par rapport aux  $n$  facteurs de Lagrange  $\alpha_i$ , qui doivent être positifs ou nuls.

Les conditions, dites de Karush-Kuhn-Tucker, sont satisfaites à l'optimum :

$$\alpha_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, n.$$

# SUPPORT VECTOR MACHINES (HARD MARGIN)

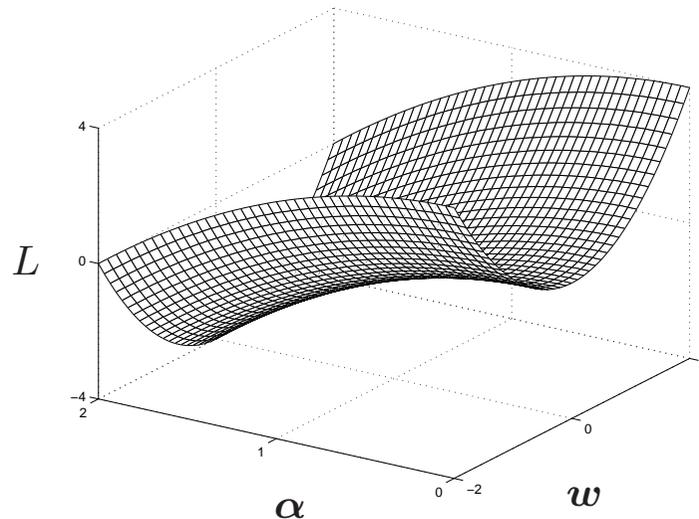
Résolution par la méthode de Lagrange

---

On résout le problème précédent à l'aide de la méthode du Lagrangien

$$L(\mathbf{w}, b; \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i \{y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1\}, \quad \alpha_i \geq 0.$$

La fonction  $L$  doit être minimisée par rapport aux variables primales  $\mathbf{w}$  et  $b$  et maximisée par rapport aux variables duales  $\alpha_i$ .



# SUPPORT VECTOR MACHINES (HARD MARGIN)

## Formulation du problème dual

---

Les conditions d'optimalité formulées à l'égard du Lagrangien

$$L(\mathbf{w}, b; \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i \{y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1\}$$

se traduisent par des dérivées nulles par rapport aux variables primales et duales :

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b; \boldsymbol{\alpha}) = 0 \quad \frac{\partial}{\partial b} L(\mathbf{w}, b; \boldsymbol{\alpha}) = 0.$$

Un rapide calcul mène aux relations suivantes qui, injectées dans l'expression du Lagrangien, fourniront le problème dual à résoudre.

$$\sum_{i=1}^n \alpha_i^* y_i = 0 \quad \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i.$$

# SUPPORT VECTOR MACHINES (HARD MARGIN)

## Formulation du problème dual

---

Le problème d'optimisation dual s'exprime finalement ainsi :

$$\text{Maximiser } W(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{sous les contraintes } \sum_{i=1}^n \alpha_i y_i = 0, \quad \alpha_i \geq 0, \quad \forall i = 1, \dots, n.$$

# SUPPORT VECTOR MACHINES (HARD MARGIN)

Formulation de la solution et support vectors

---

Le vecteur normal au plan séparateur optimum s'exprime sous la forme :

$$\mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$$

D'après les conditions de Karush-Kuhn-Tucker, à l'optimum on a :

$$\alpha_i^* \{y_i(\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) - 1\} = 0.$$

**Cas 1** :  $y_i(\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) > 1$

On a  $\alpha_i^* = 0$ , signifiant que  $\mathbf{x}_i$  n'apparaît pas dans l'expression de  $\mathbf{w}^*$ .

**Cas 2** :  $y_i(\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) = 1$

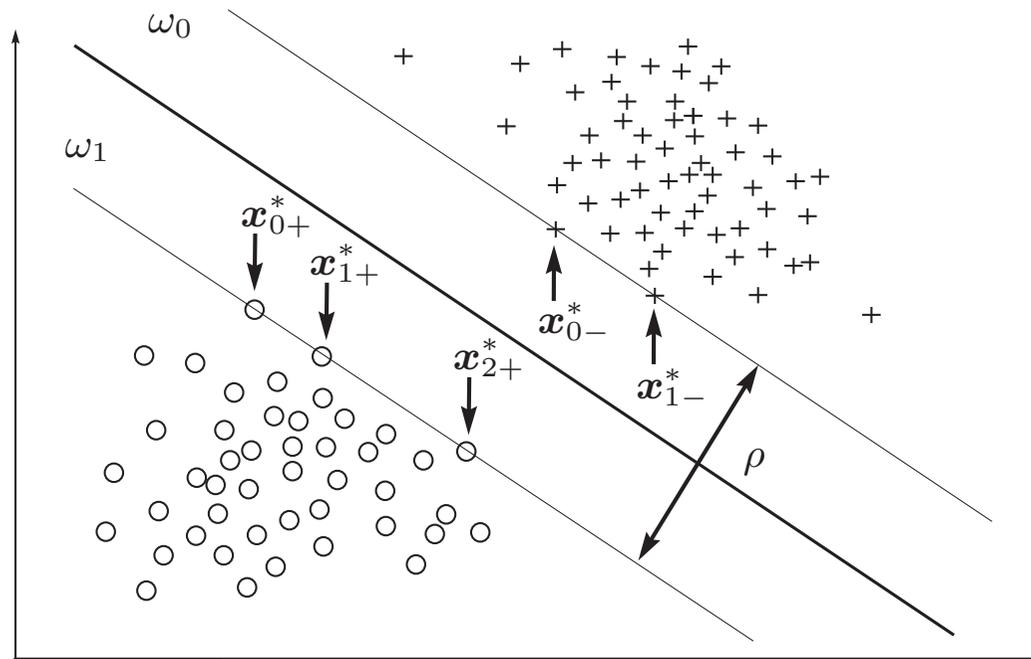
On a  $\alpha_i^* \neq 0$  et  $\mathbf{x}_i$  se trouve sur la marge. On déduit  $b^*$  de tels points.

Le vecteur  $\mathbf{w}^*$  n'est défini qu'à partir des  $\mathbf{x}_i$  situé sur la marge, les *Support Vectors*.

# SUPPORT VECTOR MACHINES (HARD MARGIN)

## Support vectors

Les support vectors sont indiqués ci-dessous par les flèches.



# PERFORMANCES EN GÉNÉRALISATION DES SVM

---

Le fait que l'hyperplan optimum s'exprime uniquement à partir des support vectors est remarquable car, en général, leur nombre est faible.

Le nombre  $n_{sv}$  de support vectors permet d'estimer les performances en généralisation du classifieur :

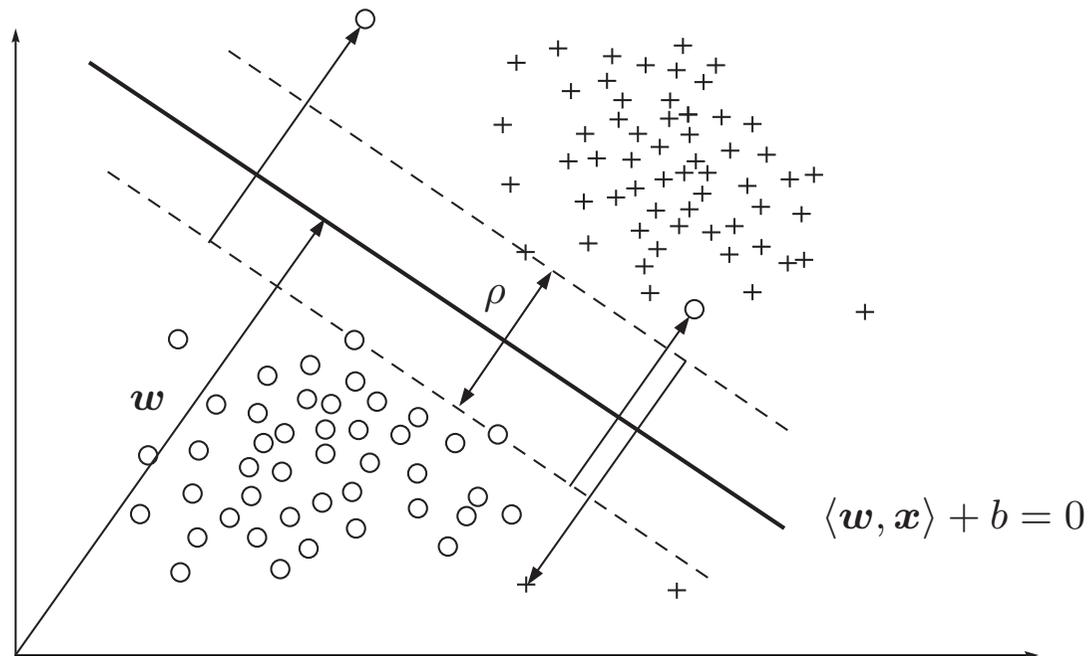
$$\mathbb{E}\{P_e\} \leq \frac{\mathbb{E}\{n_{sv}\}}{n}$$

# SUPPORT VECTOR MACHINES (SOFT MARGIN)

## Mode de pénalisation

---

Lorsque les classes en compétition ne sont pas séparables linéairement, il convient de modifier la formulation du problème afin de pénaliser les données mal classées.

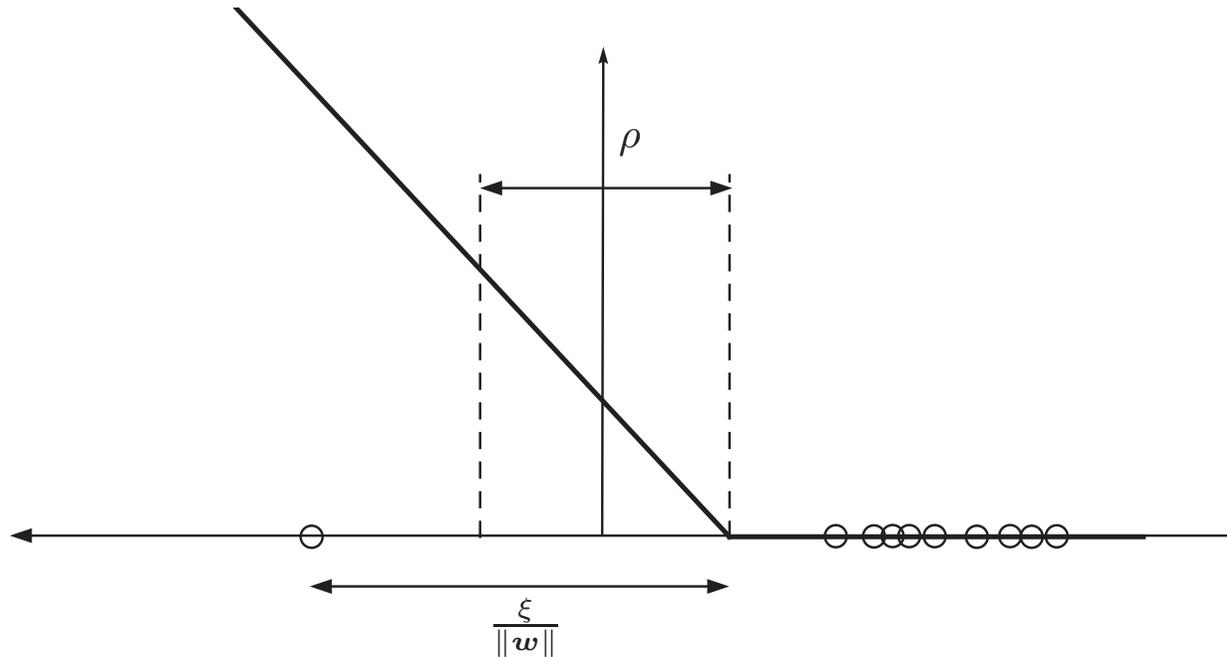


# SUPPORT VECTOR MACHINES (SOFT MARGIN)

## Fonctions de pénalisation

---

Le mode de pénalisation le plus courant est relatif à la distance de l'échantillon mal-classé à la marge. On considère parfois le carré de cette dernière.



# SUPPORT VECTOR MACHINES (SOFT MARGIN)

## Formulation du problème d'optimisation

---

Le schéma précédent incite à formuler le problème d'optimisation ainsi.

$$\text{Minimiser } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i, \quad C \geq 0$$

$$\text{sous les contraintes } y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad 1 \leq i \leq n.$$

Le terme  $C \sum_{i=1}^n \xi_i$  a pour effet de pénaliser les échantillons mal classés. D'autres fonctions de pénalisation que celle-ci existent.

# SUPPORT VECTOR MACHINES (SOFT MARGIN)

Résolution par la méthode de Lagrange

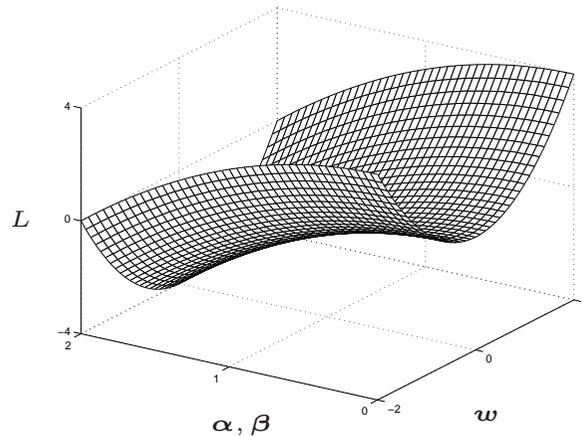
---

On résout le problème précédent à l'aide de la méthode du Lagrangien

$$L(\mathbf{w}, b, \boldsymbol{\xi}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i \{y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i\} - \sum_{i=1}^n \beta_i \xi_i,$$

où les  $\alpha_i$  et  $\beta_i$  sont des facteurs de Lagrange, positifs ou nuls.

La fonction  $L$  doit être minimisée par rapport aux variables primales  $\mathbf{w}$  et  $b$  et maximisée par rapport aux variables duales  $\alpha_i$  et  $\beta_i$ .



# SUPPORT VECTOR MACHINES (SOFT MARGIN)

## Formulation du problème dual

---

Les conditions d'optimalité formulées à l'égard du Lagrangien

$$L(\mathbf{w}, b, \boldsymbol{\xi}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i \{y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1 + \xi_i\} - \sum_{i=1}^n \beta_i \xi_i,$$

se traduisent par des dérivées nulles par rapport aux variables primales et duales :

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\xi}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = 0 \implies \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$$

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \boldsymbol{\xi}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = 0 \implies \sum_{i=1}^n \alpha_i^* y_i = 0$$

$$\frac{\partial}{\partial \boldsymbol{\xi}} L(\mathbf{w}, b, \boldsymbol{\xi}; \boldsymbol{\alpha}, \boldsymbol{\beta}) = 0 \implies \beta_i^* = C - \alpha_i^*$$

Injectées dans l'expression du Lagrangien, ces relations fournissent le problème dual à résoudre.

# SUPPORT VECTOR MACHINES (SOFT MARGIN)

## Formulation du problème dual et solution

---

Le problème d'optimisation dual s'exprime finalement ainsi :

$$\begin{aligned} \text{Minimiser } W(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{sous les contraintes } &\sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, n. \end{aligned}$$

La solution du problème s'écrit finalement

$$d(\mathbf{x}; \boldsymbol{\alpha}^*, b^*) = \text{sign} \left( \sum_{sv} \alpha_i^* y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b^* \right)$$

Afin de déterminer  $b^*$ , on utilise les conditions de Karush-Kuhn-Tucker :

$$\alpha_i^* \{ y_i (\langle \mathbf{w}^*, \mathbf{x}_i \rangle + b^*) - 1 + \xi_i^* \} = 0, \quad \beta_i^* \xi_i^* = 0.$$

Pour tout support vector  $\mathbf{x}_i$  tel que  $\alpha_i < C$ , on a  $\xi_i = 0$  et  $b^* = y_i - \langle \mathbf{w}^*, \mathbf{x}_i \rangle$ .

# SUPPORT VECTOR MACHINES (SOFT MARGIN)

## Choix du paramètre $C$

---

Minimiser  $\frac{1}{2}\|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$ ,  $C \geq 0$

sous les contraintes  $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i$ ,  $\xi_i \geq 0$ ,  $1 \leq i \leq n$ .

Le paramètre  $C$  établit un compromis entre la largeur de la marge, qui a un rôle régularisant, et le nombre d'échantillons mal classés.

**C grand** : petite marge, moins d'erreurs de classification

**C petit** : grande marge, plus d'erreurs de classification

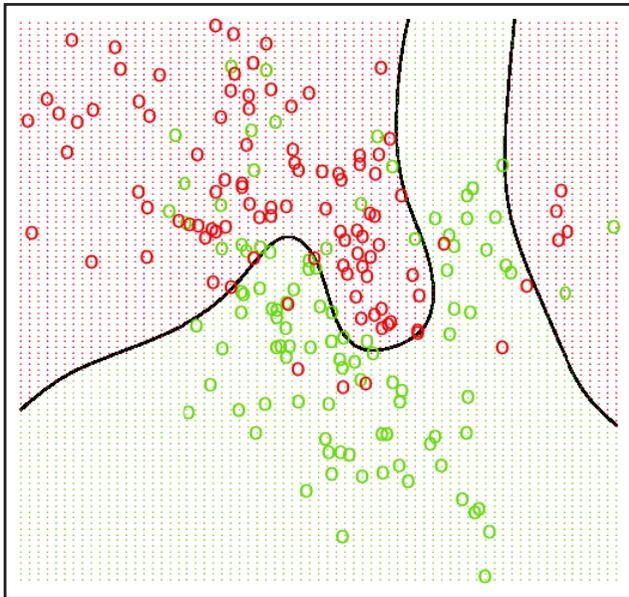
Le choix du paramètre  $C$  peut être optimisé par validation croisée.

# EXEMPLES DE MISES EN ŒUVRE

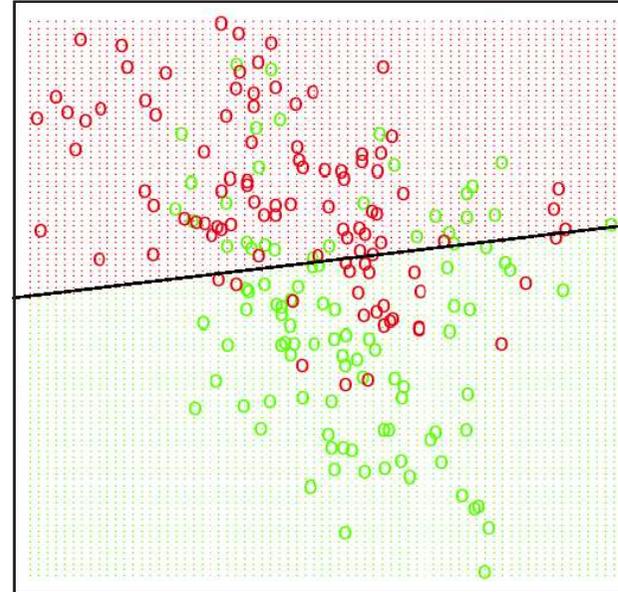
Détecteur de Bayes et discriminant de Fisher

---

Classifieur de Bayes



Régression 0/1

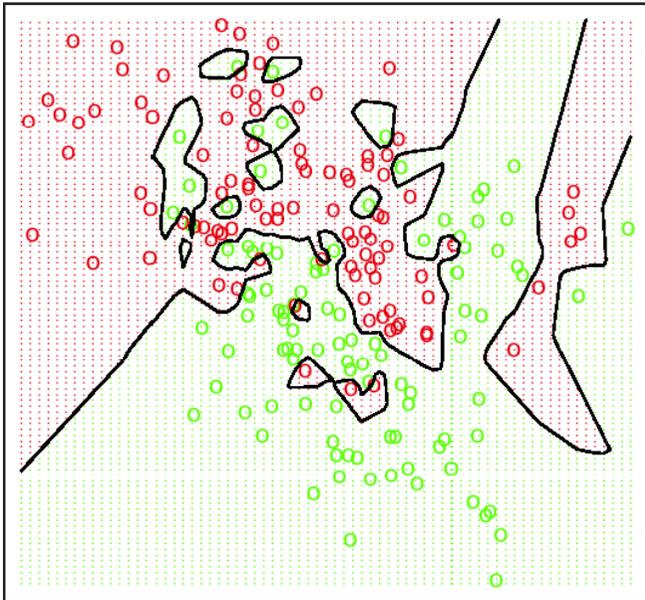


# EXEMPLES DE MISES EN ŒUVRE

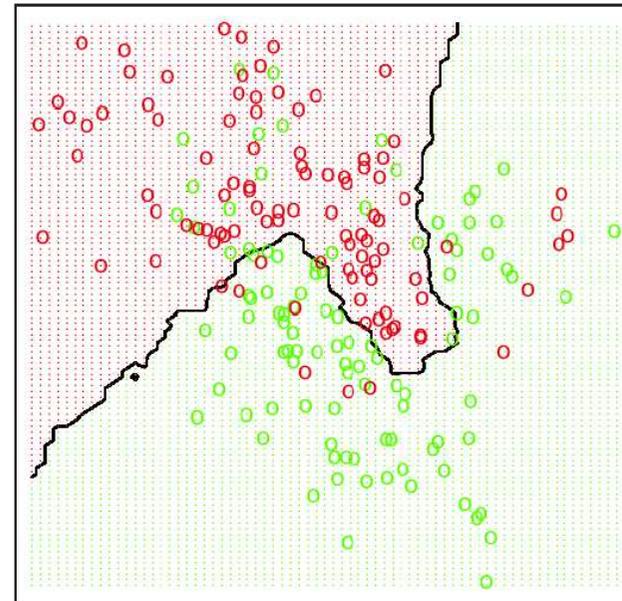
K plus-proches-voisins

---

1-ppv



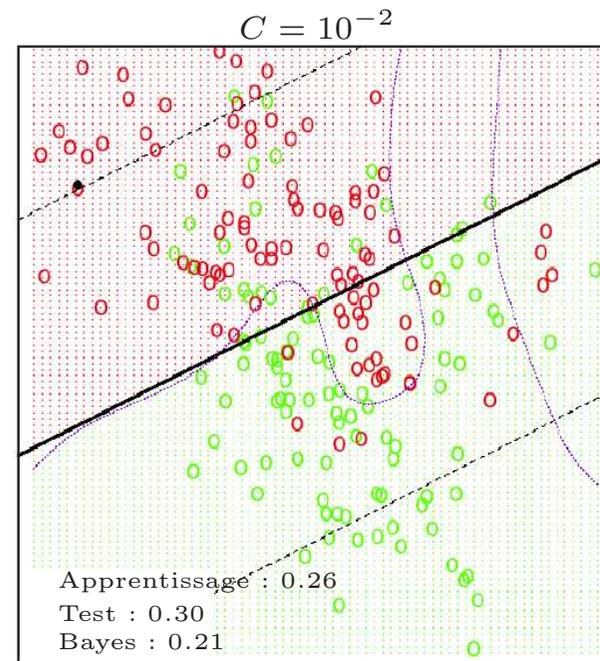
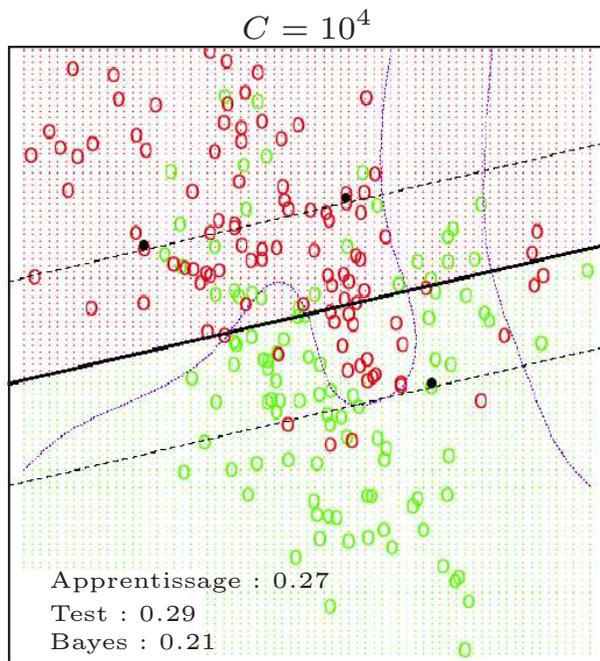
15-ppv



# EXEMPLES DE MISES EN ŒUVRE

Support vector machines

---



# SUPPORT VECTOR MACHINES

Du linéaire au non-linéaire

---

Les classifieurs linéaires ont des capacités de classification limitées. Pour y remédier, on peut les mettre en œuvre après transformation non-linéaire des données :

$$\mathbf{x} \longrightarrow \phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots]^t$$

où les  $\phi_i(\mathbf{x})$  sont des fonctions non-linéaires préalablement choisies.

**Un classifieur linéaire en  $\phi(\mathbf{x})$  est non-linéaire par rapport à  $\mathbf{x}$**

# SUPPORT VECTOR MACHINES

Exemple : classifieur polynomial

---

On considère  $\mathbf{x} = [x(1) \ x(2) \ x(3)]^t$ . Considérons la transformation suivante :

$$\begin{aligned}\phi_1(\mathbf{x}) &= x(1) & \phi_4(\mathbf{x}) &= x(1)^2 & \phi_7(\mathbf{x}) &= x(1) x(2) \\ \phi_2(\mathbf{x}) &= x(2) & \phi_5(\mathbf{x}) &= x(2)^2 & \phi_8(\mathbf{x}) &= x(1) x(3) \\ \phi_3(\mathbf{x}) &= x(3) & \phi_6(\mathbf{x}) &= x(3)^2 & \phi_9(\mathbf{x}) &= x(2) x(3)\end{aligned}$$

Un classifieur linéaire dans l'espace transformé  $\{\phi(\mathbf{x})\}_{\mathbf{x} \in \mathbb{R}^3}$ , c'est-à-dire

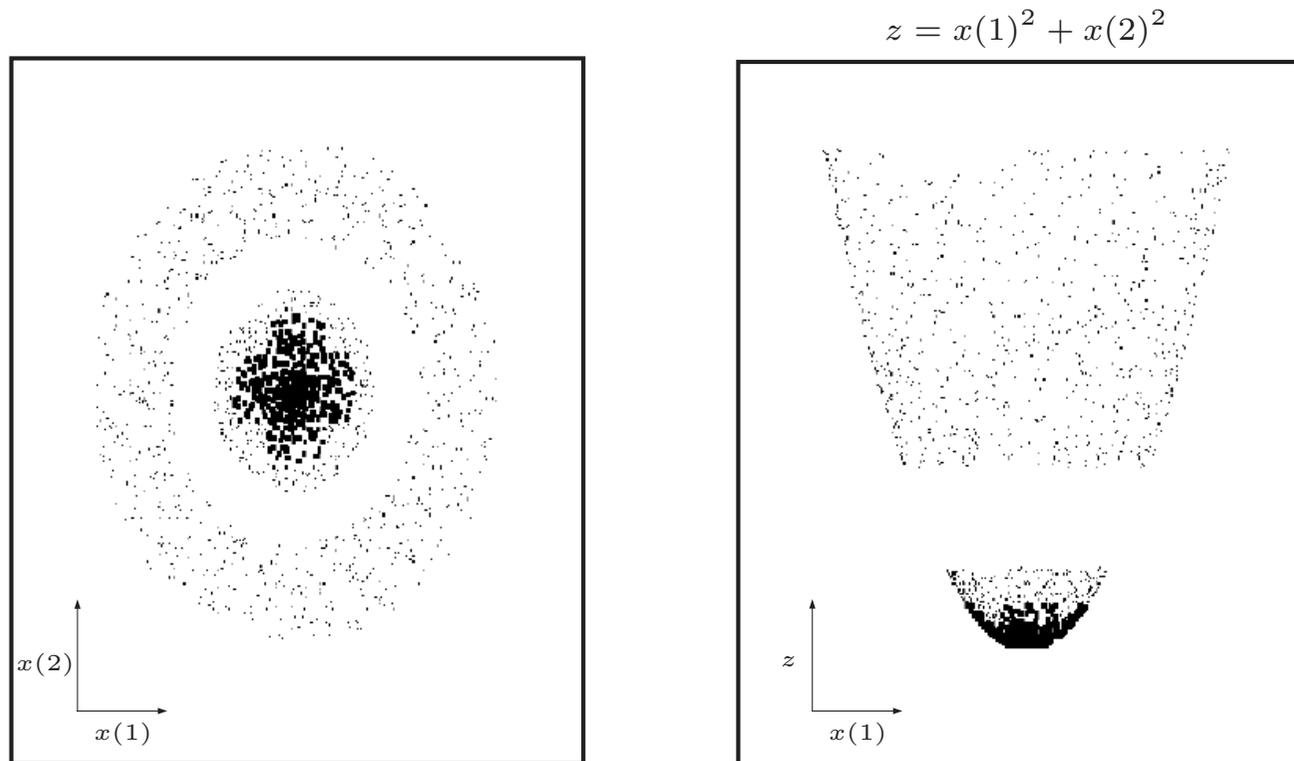
$$d(\mathbf{x}; \mathbf{w}, b) = \text{sign}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b),$$

est un classifieur polynomial de degré 2 par rapport à  $\mathbf{x}$ .

# SUPPORT VECTOR MACHINES

Exemple : classifieur polynomial

---



La transformation polynomiale rend les données linéairement séparables.

# SUPPORT VECTOR MACHINES

## Optimisation dans un espace transformé

---

Le problème d'optimisation dual s'exprime ainsi

$$\text{Minimiser } W(\boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle \boldsymbol{\phi}(\mathbf{x}_i), \boldsymbol{\phi}(\mathbf{x}_j) \rangle$$

sous les contraintes  $\sum_{i=1}^n \alpha_i y_i = 0, \quad 0 \leq \alpha_i \leq C, \quad \forall i = 1, \dots, n.$

La solution du problème s'écrit

$$d(\mathbf{x}; \boldsymbol{\alpha}^*, b^*) = \text{sign} \left( \sum_{sv} \alpha_i^* y_i \langle \boldsymbol{\phi}(\mathbf{x}), \boldsymbol{\phi}(\mathbf{x}_i) \rangle + b^* \right).$$

On remarque que

- on n'a jamais besoin de calculer explicitement  $\boldsymbol{\phi}(\mathbf{x})$  ;
- si  $\mathbf{x}$  est de grande dimension, celle de  $\boldsymbol{\phi}(\mathbf{x})$  l'est encore d'avantage, parfois infinie.

# SUPPORT VECTOR MACHINES

## Le coup du noyau

---

Si l'on parvient à définir un noyau  $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$  tel que :

— la surface de décision associée soit performante

$$d(\mathbf{x}; \boldsymbol{\alpha}^*, b^*) = \text{sign} \left( \sum_{sv} \alpha_i^* y_i \kappa(\mathbf{x}_i, \mathbf{x}_j) + b^* \right)$$

— il soit aisé de calculer  $\kappa(\mathbf{x}_i, \mathbf{x}_j)$ , même pour des données de grande dimension...

**Le tour est joué !**

# LE COUP DU NOYAU

## Noyaux polynomiaux

---

Dans le cas de la transformation polynomiale de degré 2, on montre aisément que :

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = (1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^2 \triangleq \kappa(\mathbf{x}, \mathbf{x}')$$

▷ Le calcul de produit scalaire peut s'effectuer dans  $\mathbb{R}^2$  !

Plus généralement, on s'intéresse à  $\kappa(\mathbf{x}, \mathbf{x}') = (1 + \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle)^q$ , avec  $\mathbf{x} \in \mathbb{R}^l$ .

$$\kappa(\mathbf{x}, \mathbf{x}') = (1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^q = \sum_{j=0}^q \binom{q}{j} \langle \mathbf{x}, \mathbf{x}' \rangle^j.$$

Chaque composante  $\langle \mathbf{x}, \mathbf{x}' \rangle^j = [x(1)x'(1) + \dots + x(l)x'(l)]^j$  de cette expression peut être développée en une somme pondérée de monômes de degré  $j$  de la forme

$$[x(1)x'(1)]^{j_1} [x(2)x'(2)]^{j_2} \dots [x(l)x'(l)]^{j_l}$$

avec  $\sum_{i=1}^l j_i = j$ . Ceci mène directement à l'expression de  $\phi(\mathbf{x})$ ...

# LE COUP DU NOYAU

## Théorème de Mercer

---

On s'intéresse aux fonctions  $\kappa(\mathbf{x}, \mathbf{x}')$  pouvant faire fonction de produit scalaire dans un espace  $\mathcal{H}$ . On appelle *noyau* une fonction symétrique  $\kappa$  de  $\mathcal{X} \times \mathcal{X}$  dans  $\mathbb{R}$ .

**Théorème 2.** *Si  $\kappa$  est un noyau continu d'un opérateur intégral défini positif, ce qui signifie que*

$$\iint \varphi(\mathbf{x}) \kappa(\mathbf{x}, \mathbf{x}') \varphi^*(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \geq 0$$

*pour tout  $\varphi \in \mathcal{L}^2(\mathcal{X})$ , il peut être décomposé sous la forme*

$$\kappa(\mathbf{x}, \mathbf{x}') = \sum_{i=1}^{\infty} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{x}'),$$

*où  $\psi_i$  et  $\lambda_i$  sont les fonctions propres (orthogonales) et valeurs propres (positives) du noyau  $\kappa$ , respectivement, telles que*

$$\int \kappa(\mathbf{x}, \mathbf{x}') \psi_i(\mathbf{x}) d\mathbf{x} = \lambda_i \psi_i(\mathbf{x}').$$

# LE COUP DU NOYAU

## Théorème de Mercer

---

Il est aisé de voir qu'un noyau  $\kappa$  satisfaisant au théorème de Mercer peut faire fonction de produit scalaire dans un espace transformé  $\mathcal{H}$ . Il suffit d'écrire :

$$\phi(\mathbf{x}) = \begin{pmatrix} \sqrt{\lambda_1} \psi_1(\mathbf{x}) \\ \sqrt{\lambda_2} \psi_2(\mathbf{x}) \\ \dots \end{pmatrix}$$

Dans ces conditions, on vérifie bien que l'on retrouve :  $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle = \kappa(\mathbf{x}, \mathbf{x}')$ .

On définit l'espace  $\mathcal{H}$  comme étant celui engendré par les fonctions propres  $\psi_i$  du noyau  $\kappa$ , c'est-à-dire

$$\mathcal{H} \triangleq \{f(\cdot) \mid f(x) = \sum_{i=1}^{\infty} \alpha_i \psi_i(x), \alpha_i \in \mathbb{R}\}.$$

# LE COUP DU NOYAU

## Exemples de noyaux de Mercer

---

On peut montrer que les noyaux suivants vérifie la condition de Mercer, et correspondent donc à un produit scalaire dans un espace  $\mathcal{H}$ .

Noyaux projectifs	
sigmoïdal	$\frac{1}{\eta_0} \tanh(\beta_0 \langle \mathbf{x}, \mathbf{x}' \rangle - \alpha_0)$
monomial de degré $q$	$\langle \mathbf{x}, \mathbf{x}' \rangle^q$
polynomial de degré $q$	$(1 + \langle \mathbf{x}, \mathbf{x}' \rangle)^q$

Noyaux radiaux	
uniforme	$\frac{1}{\eta_0} \mathbb{1}_{\ \mathbf{x} - \mathbf{x}'\  \leq \beta_0}$
Epanechnikov	$\frac{1}{\eta_0} (\beta_0^2 - \ \mathbf{x} - \mathbf{x}'\ ^2) \mathbb{1}_{\ \mathbf{x} - \mathbf{x}'\  \leq \beta_0}$
Cauchy	$\frac{1}{\eta_0} \frac{1}{1 + \ \mathbf{x} - \mathbf{x}'\ ^2 / \beta_0^2}$

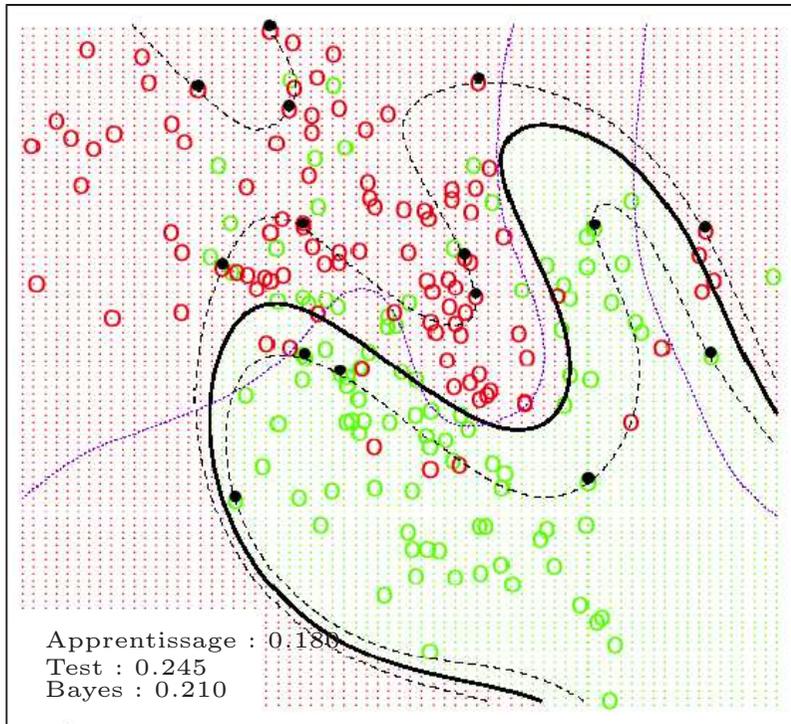
... et encore  $\kappa_1(\mathbf{x}, \mathbf{x}') + \kappa_2(\mathbf{x}, \mathbf{x}')$ ,  $\kappa_1(\mathbf{x}, \mathbf{x}') \cdot \kappa_2(\mathbf{x}, \mathbf{x}')$ , ...

# EXEMPLES DE MISES EN ŒUVRE

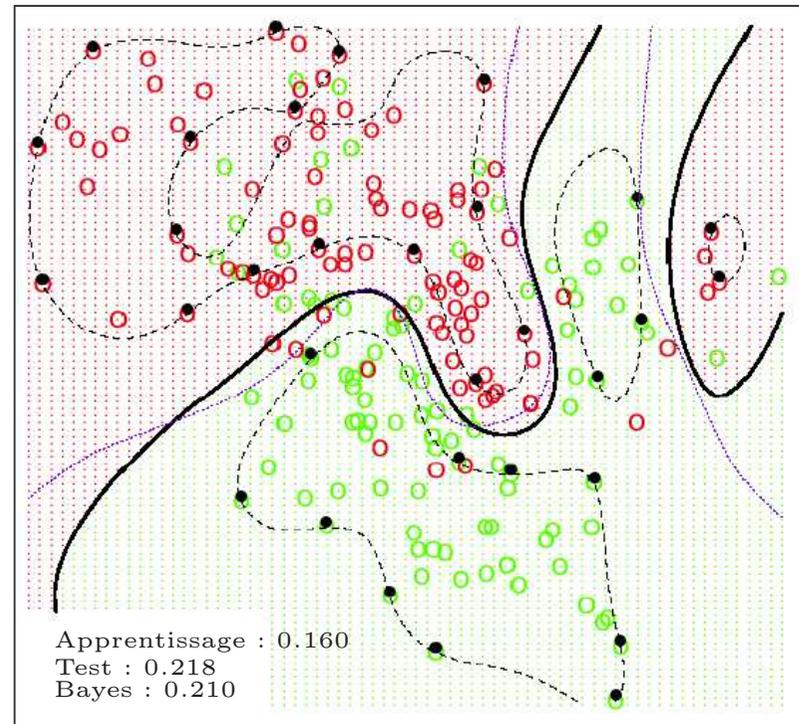
Support vector machines

---

noyau polynomial de degré 4



noyau gaussien



# LE COUP DU NOYAU

En résumé

---

Parmi les possibilités offertes par le coup du noyau, on note en particulier :

- Tout algorithme pouvant s'exprimer à partir de produits scalaires peut bénéficier du coup du noyau en vue d'une extension au cas non-linéaire.
- Grâce au coup du noyau, tout algorithme de reconnaissance des formes est en mesure de traiter des données autres que numériques, par exemple alphabétiques.

# PREMIER BILAN SUR LES SVM

## Qualités

---

Par rapport aux techniques concurrentes telles que les réseaux de neurones artificiels, les SVM possèdent d'immenses qualités.

1. Solution unique  
→ Problème quadratique
2. Processus de régularisation intégré, solution sparse  
→ Fonction coût et contraintes d'inégalité en découlant
3. Extension aisée au cas non-linéaire, solution non boîte noire  
→ Kernel trick

# PREMIER BILAN SUR LES SVM

Une difficulté majeure

---

La mise en œuvre de l'algorithme d'apprentissage par une approche directe s'avère difficile car la taille du problème est celle de la base d'apprentissage  $\mathcal{A}_n$ .

Minimiser  $W(\boldsymbol{\alpha}) = \frac{1}{2} \boldsymbol{\alpha}^t \mathbf{H} \boldsymbol{\alpha} + \mathbf{1}_n^t \boldsymbol{\alpha}$

sous les contraintes  $\mathbf{y}^t \boldsymbol{\alpha} = 0, \quad 0 \cdot \mathbf{1}_n \leq \boldsymbol{\alpha} \leq C \cdot \mathbf{1}_n.$

# DES SOLUTIONS ALGORITHMIQUES

## Méthodes par décomposition

---

Des algorithmes performants existent, basés sur la décomposition du problème d'optimisation en sous-problèmes.

Minimiser

$$W(\alpha_B) = \frac{1}{2} (\alpha_B \ \alpha_N)^t \begin{pmatrix} \mathbf{H}_{BB} & \mathbf{H}_{BN} \\ \mathbf{H}_{NB} & \mathbf{H}_{NN} \end{pmatrix} (\alpha_B \ \alpha_N) - (\mathbf{1}_B \ \mathbf{1}_N) \begin{pmatrix} \alpha_B \\ \alpha_N \end{pmatrix}$$

sous les contraintes  $\alpha_B^t \mathbf{y}_B + \alpha_N^t \mathbf{y}_N = 0 \quad 0 \cdot \mathbf{1}_B \leq \alpha_B \leq C \cdot \mathbf{1}_B.$

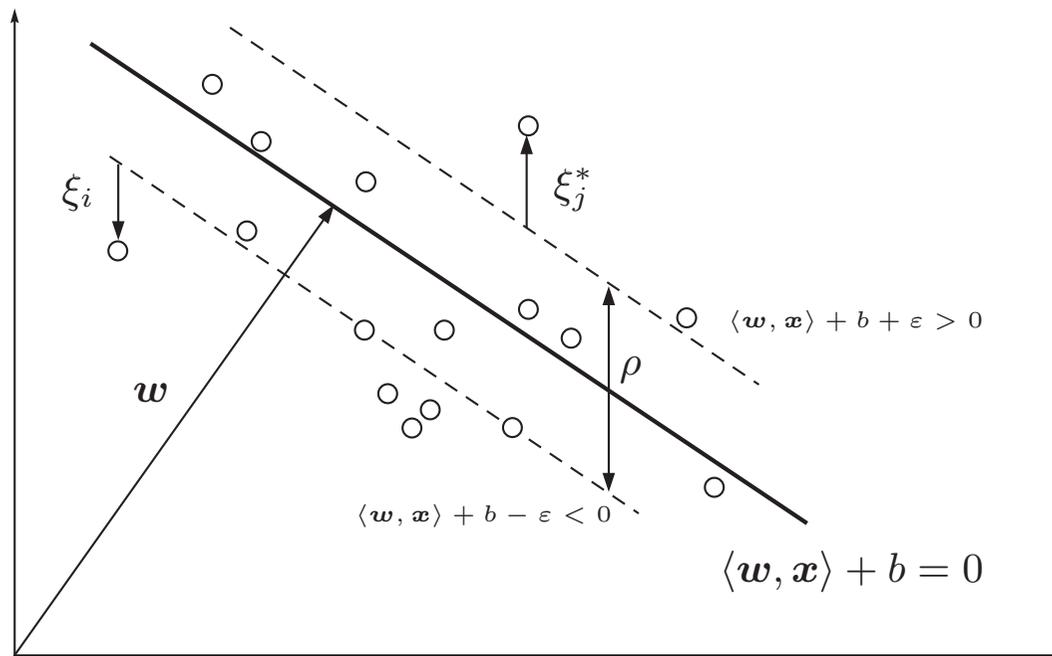
Plusieurs stratégies ont été proposées pour décomposer le problème. La plus extrême consiste à limiter  $B$  à deux éléments. La résolution est alors analytique.

# APPLICATION À LA RÉGRESSION

## Support vector regression

---

L'idée de marge, sur laquelle repose la presque totalité des qualités des SVM, peut être appliquée à la résolution de problèmes de régression.



# APPLICATION À LA RÉGRESSION

## Support vector regression

---

Le problème d'optimisation correspondant s'exprime de la façon suivante.

$$\text{Minimiser } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*)$$

$$\text{sous les contraintes } y_i - (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \varepsilon + \xi_i, \quad \xi_i \geq 0$$

$$(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - y_i \geq \varepsilon + \xi_i^*, \quad \xi_i^* \geq 0$$

Le kernel trick s'applique de la même façon que pour les SVM, tout comme les algorithmes de résolution.

# APPLICATION À LA RÉGRESSION

## Support vector regression

---

On résout le problème précédent à l'aide de la méthode du Lagrangien. On aboutit au problème dual suivant :

$$\begin{aligned} \text{Minimiser } W(\boldsymbol{\alpha}, \boldsymbol{\alpha}^*) &= \frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) + \langle \boldsymbol{x}_i, \boldsymbol{x}_j \rangle \\ &+ \varepsilon \sum_{i=1}^n (\alpha_i + \alpha_i^*) + \sum_{i=1}^n y_i (\alpha_i - \alpha_i^*) \end{aligned}$$

sous les contraintes

$$\sum_{i=1}^n \alpha_i = \sum_{i=1}^n \alpha_i^*, \quad 0 \leq \alpha_i, \alpha_i^* \leq C, \quad \forall i = 1, \dots, n.$$

La solution s'écrit finalement sous la forme suivante, permettant ainsi la mise en œuvre du coup du noyau :

$$f(\boldsymbol{x}) = \sum_{i=1}^n (\alpha_i^* - \alpha_i) \langle \boldsymbol{x}, \boldsymbol{x}_i \rangle.$$

# APPLICATION À LA RÉGRESSION

## Support vector regression

