

# Une introduction aux machines à vecteurs supports (SVM)

# Plan

---

- ▶ Historique
- ▶ Quelle est la bonne frontière de séparation pour deux classes linéairement séparables ?
  - ▶ La solution SVM
- ▶ Adaptation aux cas non linéairement séparables: l'astuce des fonctions noyau
- ▶ Exemples d'application
- ▶ Conclusion

# Historique du SVM

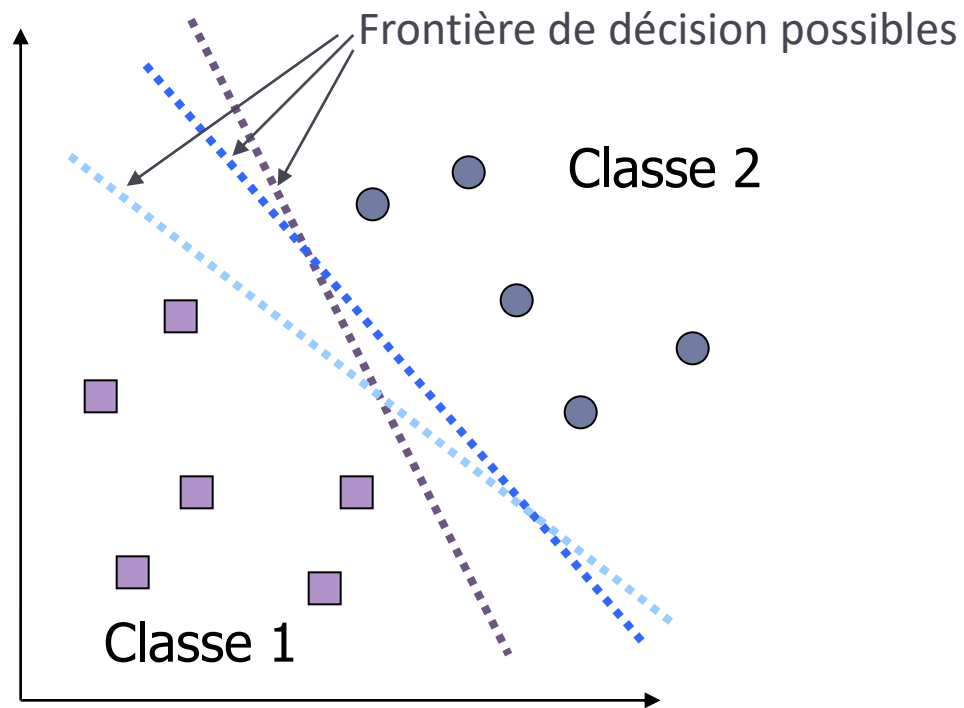
---

- ▶ Classifieur dérivé de la théorie de l'apprentissage statistique par Vapnik et Chervonenkis (~1994)
- ▶ Devenu populaire depuis qu'il a permis des performances égales ou meilleures aux RNA pour reconnaître l'écriture manuscrite en partant d'images formées de pixels.
- ▶ Proche de :
  - ▶ Séparateurs à vastes marges
  - ▶ Méthodes à fonctions noyau
  - ▶ Réseaux de neurones à bases radiales

# Problème à deux classes linéairement séparables

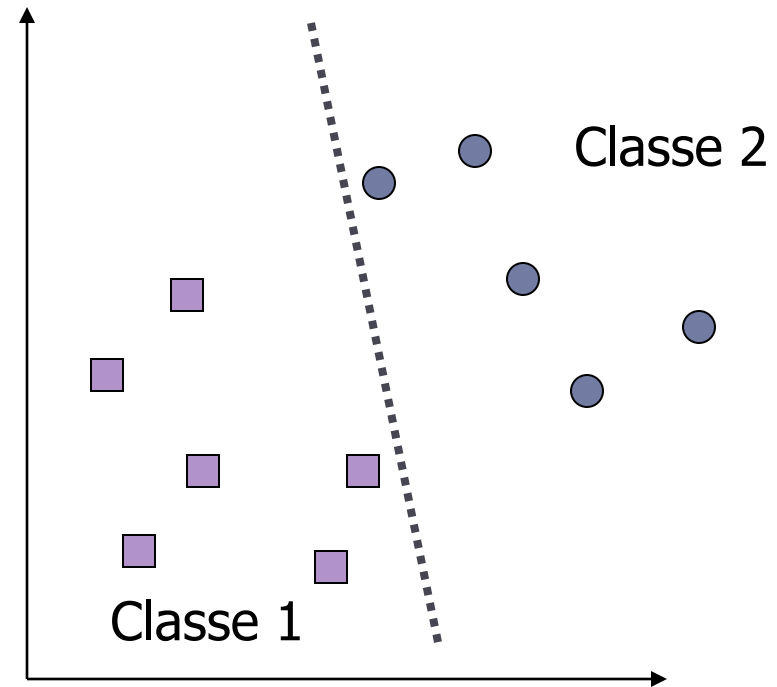
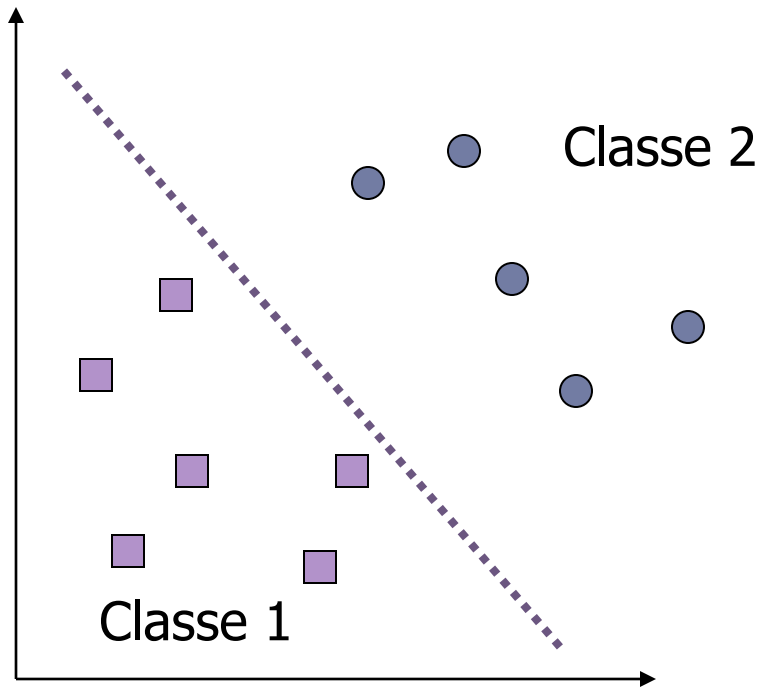
---

- ▶ Plusieurs surfaces de décision existent pour séparer les classes ; laquelle choisir ?



# Exemples de mauvais choix

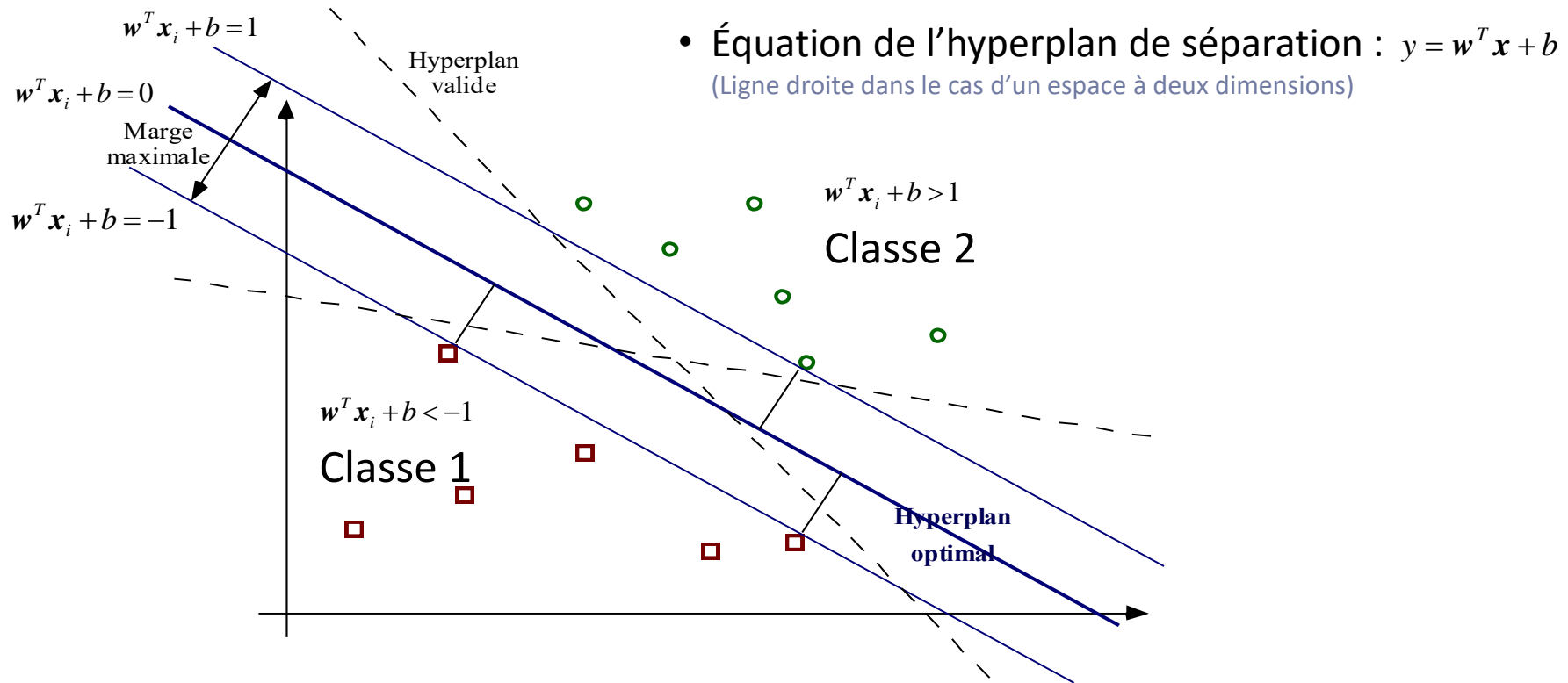
---



- Pour minimiser la sensibilité au bruit, la droite de décision doit être aussi éloignée que possible des données de chaque classe



# Hyperplan de plus vaste marge



- Si  $\{x_1, \dots, x_n\}$  est l'ensemble des données et  $y_i \in \{1, -1\}$  est la classe de chacune, on devrait avoir :

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad \forall i$$

tout en ayant une distance optimale entre chaque  $x_i$  et le plan de séparation

# Optimisation de la marge

- ▶ Distance normale d'un point à l'hyperplan :

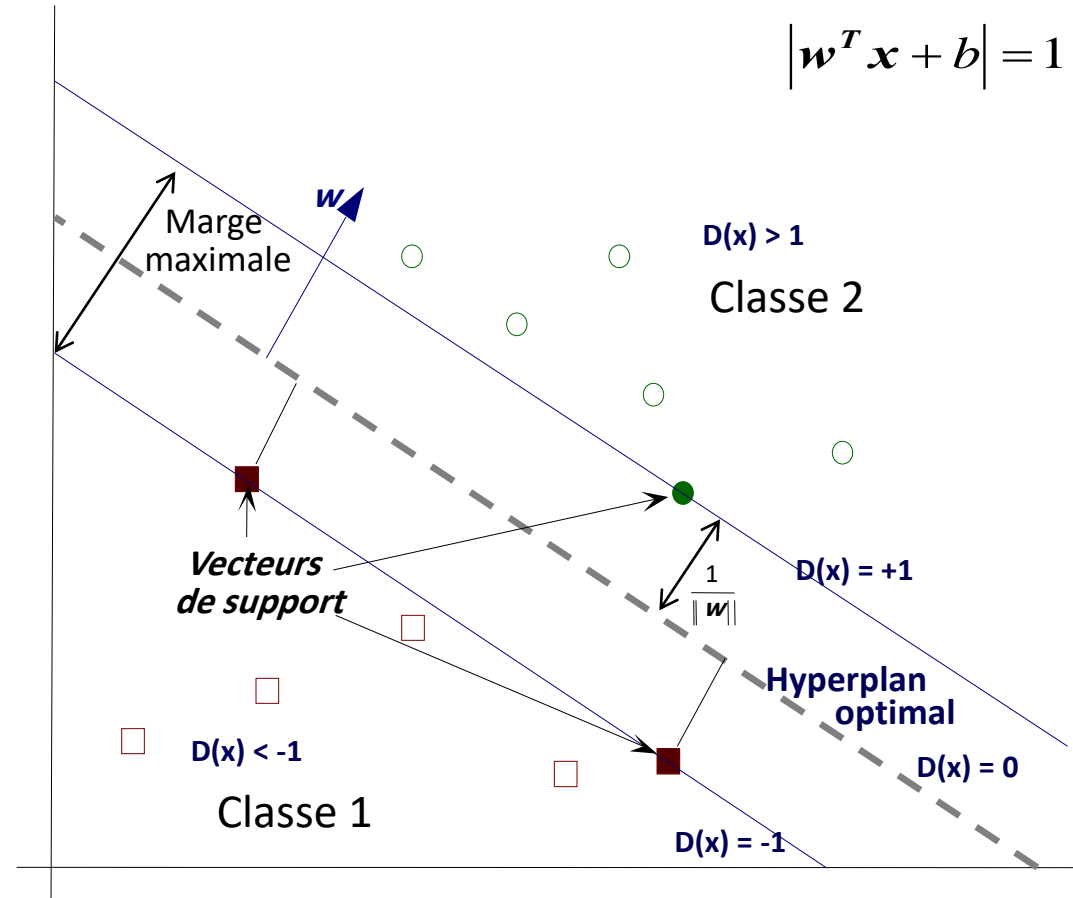
$$D(\mathbf{x}) = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|}$$

- ▶ Marge max. avant d'atteindre les frontières des deux classes ( $|\mathbf{w}^T \mathbf{x} + b| = 1$ ) :

$$m = \frac{2}{\|\mathbf{w}\|}$$

- ▶ Maximiser  $m$  revient à minimiser  $\|\mathbf{w}\|$  tout en préservant le pouvoir de classification :

$$\min \frac{1}{2} \|\mathbf{w}\|^2 \text{ sous la contrainte: } y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \forall i$$



# Problème d'optimisation quadratique

---

- ▶ Maximiser le pouvoir de généralisation du classeur revient à trouver  $\mathbf{w}$  et  $b$  tels que :

$$\frac{1}{2} \|\mathbf{w}\|^2 \text{ est minimum}$$

et

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i=1, \dots, n$$

- ▶ Si  $d$  est la dimension des  $\mathbf{x}_i$  (nombre d'entrées), cela revient à régler  $d+1$  paramètres (les éléments de  $\mathbf{w}$ , plus  $b$ )
  - ▶ Possible par les méthodes d'optimisation classiques (e.g., optimisation quadratique) seulement si  $d$  pas trop grand ( $< \text{qqs } 10^3$ )
  - ▶ L'approche SVM utilise les multiplicateurs de Lagrange pour une solution plus simple



# Les multiplicateurs de Lagrange en 30 s

---

- ▶ On veut maximiser ou minimiser  $f(x)$  sous la contrainte  $g(x)=0$
- ▶ Solutions possibles :
  - ▶ Résoudre  $g(x)=0$  et substituer la/les racines trouvées dans  $f(x)$  ; résoudre alors  $f'(x) = 0$  : pas toujours facile !
  - ▶ Considérer que  $f(x)$  et  $g(x)=0$  évoluent pareillement au voisinage de l'extrémum recherché. Leurs tangentes sont alors colinéaires et :
$$f'(x) = \alpha g'(x) \text{ (ou } f'(x) - \alpha g'(x) = 0), \alpha \text{ étant à déterminer}$$
  - ▶ La méthode des multiplicateurs de Lagrange regroupe la fonction à optimiser et la contrainte en une seule fonction  $\Lambda(x, \alpha) = f(x) - \alpha g(x)$ 
    - La solution de  $d\Lambda(x, \alpha)/dx = 0$  donne le point où les tangentes sont colinéaires; en même temps,  $d\Lambda(x, \alpha)/d\alpha = 0$  répond à la contrainte.



# Cas à plusieurs dimensions

---

- ▶ On veut trouver le minimum (ou maximum) d'une fonction  $f(\mathbf{x})$  en respectant un nombre de contraintes  $g_i(\mathbf{x})=0, i=1, \dots, n$
  - ▶ On peut montrer qu'à l'extrémum recherché :  $\nabla f(\mathbf{x}) = \sum_{i=1}^n \alpha_i \nabla g_i(\mathbf{x})$ ,  
(colinéarité des tangentes)
- ou encore  $\nabla f(\mathbf{x}) - \sum_{i=1}^n \alpha_i \nabla g_i(\mathbf{x}) = 0$ , pour des coefficients  $\alpha_i$  à déterminer

- ▶ Si on forme la fonction (lagrangien) :  $\Lambda(\mathbf{x}, \alpha) = f(\mathbf{x}) - \sum_{i=1}^l \alpha_i g_i(\mathbf{x})$

Alors :  $\nabla_{\mathbf{x}} \Lambda(\mathbf{x}, \alpha) = \nabla f(\mathbf{x}) - \sum_{i=1}^l \lambda_i \nabla g_i(\mathbf{x})$

et  $\nabla_{\alpha_i} \Lambda(\mathbf{x}, \alpha) = g_i(\mathbf{x})$

=> la solution de  $\nabla_{\mathbf{x}} \Lambda(\mathbf{x}, \alpha) = 0$  mène à un extrémum qui respecte les contraintes.



# Forme duale du problème d'optimisation

- ▶ Ici, la fonction à optimiser sous contraintes est celle de l'hyperplan de séparation à marge maximale, qui est défini par les paramètres  $(\mathbf{w}, b)$  correspondants.
- ▶ Donc, partant d'un ensemble de données  $\{(\mathbf{x}_i, y_i)\}$ , de l'ensemble de contraintes  $\{(\mathbf{x}_i^T \mathbf{w} + b)y_i - 1 = 0\}$  et des paramètres à optimiser  $(\mathbf{w}, b)$ , on a:

$$\begin{cases} \Lambda(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i \{(\mathbf{x}_i^T \mathbf{w} + b)y_i - 1\} \\ \forall i \quad \alpha_i \geq 0 \end{cases}$$

Le minimum recherché est donné par la solution de  $\nabla_{\mathbf{w}, b} \Lambda(\mathbf{w}, b, \alpha) = 0$

- ▶ Il existe une formulation duale du problème plus facile à résoudre :

- Théorème de Kuhn-Tucker :  $\min_{\alpha} \left( \max_{\mathbf{x}} \Lambda(\mathbf{x}, \alpha) \right) = \max_{\mathbf{x}} \left( \min_{\alpha} \Lambda(\mathbf{x}, \alpha) \right)$

⇒ Solution alternative pour trouver  $\mathbf{w}$  et  $b$  :

résoudre  $\nabla_{\alpha} \Lambda(\mathbf{w}, b, \alpha) = 0$  sous la contrainte  $\nabla_{\mathbf{w}, b} \Lambda(\mathbf{w}, b, \alpha) = 0$

# Forme duale du problème d'optimisation

---

- ▶ On résout l'équation du lagrangien par rapport à  $\alpha$  au lieu de  $\mathbf{w}, b$
- ▶ Avantage de résoudre  $\nabla_{\alpha}\Lambda(\mathbf{w}, b, \alpha) = 0$  au lieu de  $\nabla_{\mathbf{w}, b}\Lambda(\mathbf{w}, b, \alpha) = 0$  :
  - ▶ La complexité du problème d'optimisation devient proportionnelle à  $n$  (nombre de paires d'apprentissage  $(\mathbf{x}_i, y_i)$ ) et non  $d$  (la dimension de  $\mathbf{w}$  qui est donnée par celle des  $\mathbf{x}_i$ )
  - ▶ Possible d'obtenir des solutions pour des problèmes impliquant  $\approx 10^5$  exemples
  - ▶ Problème pour lequel le maximum global des  $\alpha_i$  peut toujours être trouvé

# Formulation du problème dual

▶ Partant de 
$$\begin{cases} \Lambda(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \alpha_i \{(\mathbf{x}_i^T \mathbf{w} + b) y_i - 1\} \\ \forall i \quad \alpha_i \geq 0 \end{cases}$$

$\nabla_{\mathbf{w}, b} \Lambda(\mathbf{w}, b, \boldsymbol{\alpha}) = 0$  donne 
$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

et 
$$\sum_{i=1}^n \alpha_i y_i = 0$$

▶ On a par substitution dans  $\Lambda(\mathbf{w}, b, \boldsymbol{\alpha})$  :

$$\begin{cases} \Lambda(\mathbf{w}, b, \boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \forall i \quad \alpha_i \geq 0 \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

▶ Il faut trouver  $\boldsymbol{\alpha}$  qui maximise  $\Lambda(\mathbf{w}, b, \boldsymbol{\alpha})$  : résoudre  $\nabla_{\boldsymbol{\alpha}} \Lambda(\mathbf{w}, b, \boldsymbol{\alpha}) = 0$

➤  $n$  équations linéaires homogènes à  $n$  inconnues (les composants de  $\boldsymbol{\alpha}$ )

# Solution du problème d'optimisation

$$\left\{ \begin{array}{l} \hat{\mathbf{w}} = \sum_{i=1}^{n_s} \hat{\alpha}_i y_i \mathbf{x}_i \\ \hat{b} = y_s - \sum_{i=1}^{n_s} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}_s) \\ D(\mathbf{x}) = (\hat{\mathbf{w}}^T \mathbf{x} + \hat{b}) \end{array} \right.$$

- $\hat{\cdot}$  : estimé
- $n_s$  : nombre de vecteurs de support ( $\mathbf{x}_i$  avec  $\alpha \neq 0$ )
- $(\mathbf{x}_s, y_s)$  : vecteur de support arbitraire (pour trouver  $\hat{b}$ )

- ▶ Les données  $\mathbf{x}_i$  avec  $\alpha \neq 0$  sont appelées vecteurs de support. Ils correspondent aux points les plus proches de la surface de séparation
- Dans l'expression du Lagrangien pour déterminer  $\hat{\mathbf{w}}$  et  $\hat{b}$ , les données  $\mathbf{x}$  apparaissent uniquement pour former des produits scalaires

$$\Lambda(\mathbf{w}, b, \boldsymbol{\alpha}) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

# Caractéristiques de la solution

---

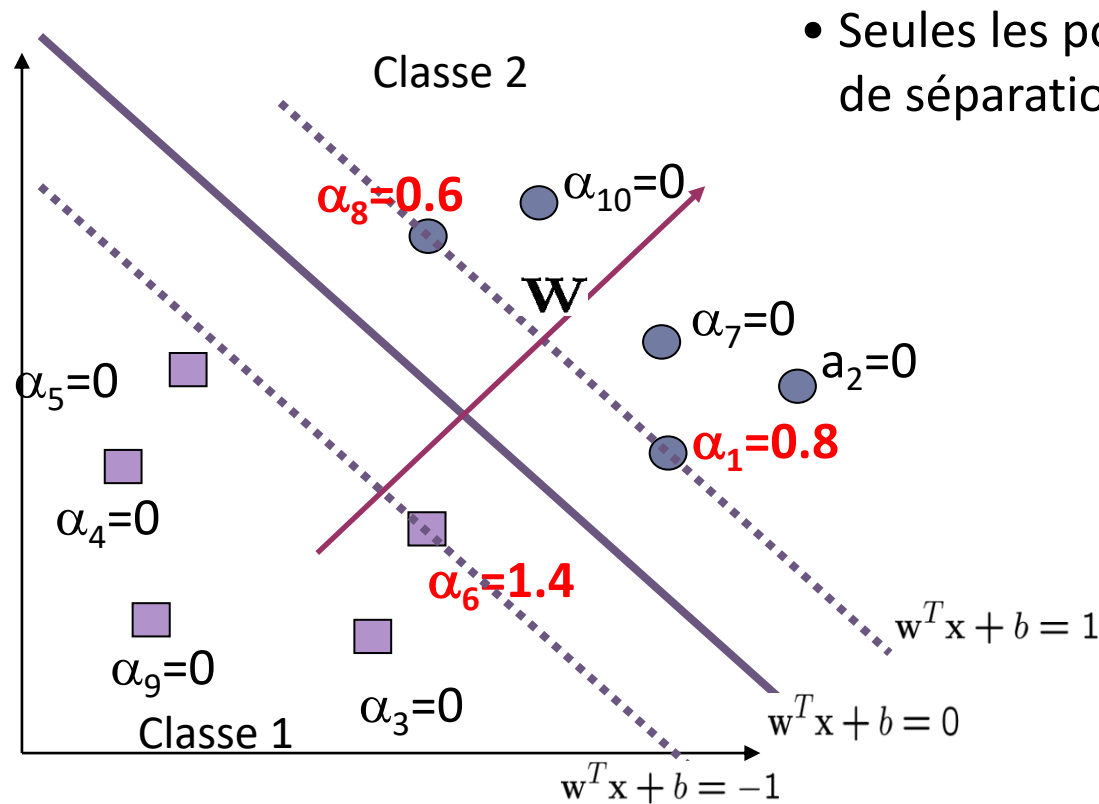
- ▶ Puisque plusieurs  $\alpha_i$  sont nuls,  $\mathbf{w}$  est une combinaison linéaire d'un petit nombre de données
- ▶ La surface de décision est uniquement déterminée par les  $n_s$  vecteurs de support trouvés:

$$\hat{\mathbf{w}} = \sum_{i=1}^{n_s} \hat{\alpha}_i y_i \mathbf{x}_i$$

$$\hat{b} = y_s - \sum_{i=1}^{n_s} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}_s)$$

- ▶ Pour classer une nouvelle donnée  $\mathbf{z}$ 
  - ▶ Calculer  $\hat{\mathbf{w}}^T \mathbf{z} + \hat{b} = \sum_{i=1}^{n_s} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{z}) + \hat{b}$  et classer  $\mathbf{z}$  dans la classe 1 si le résultat est positif et dans la classe 2 s'il est négatif

# Interprétation géométrique



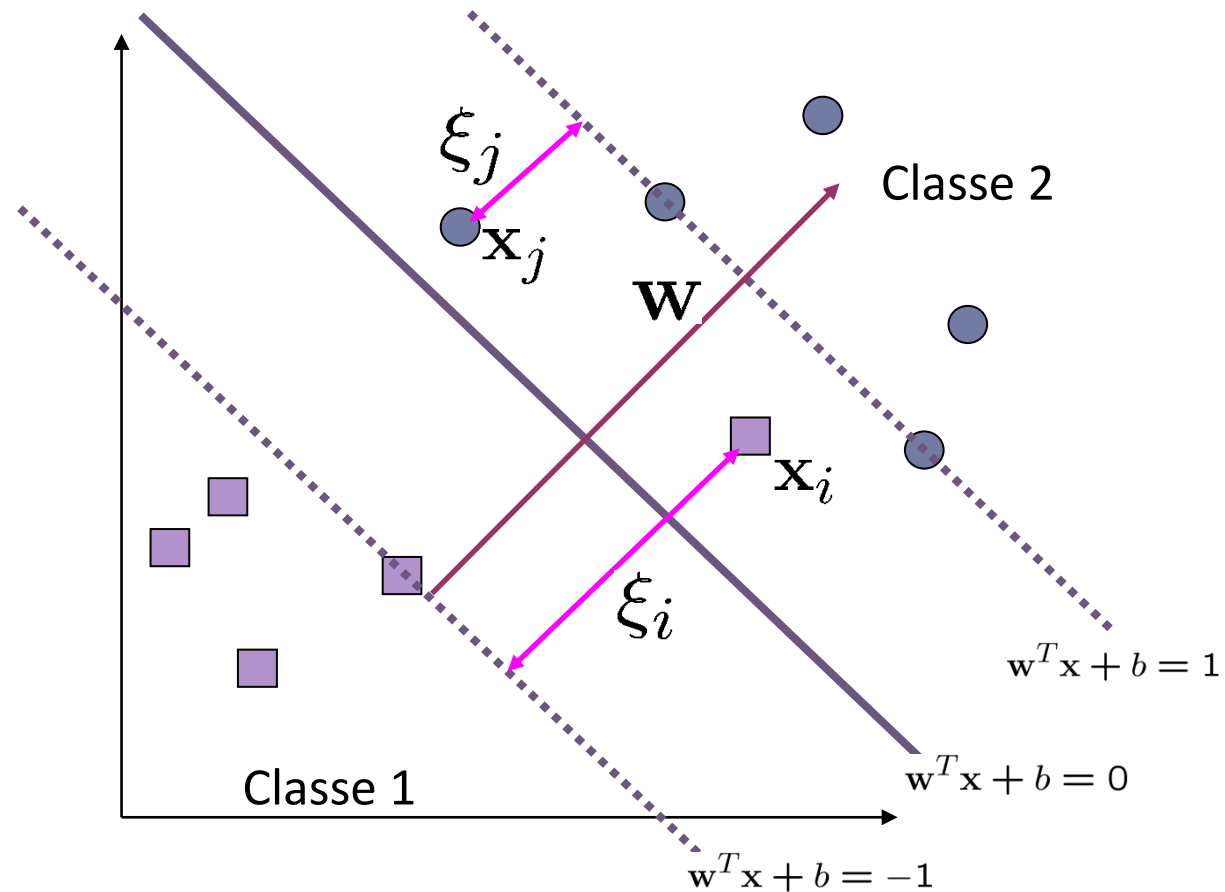
- Seules les points les plus proches de la surface de séparation influent sur sa définition

- Il existe des limites théoriques pour l'erreur de classification de données nouvelles
  - Plus grande la marge, plus petite la limite
  - Plus petit le nombre de SV, plus petite la limite



# Et pour un cas non linéairement séparable ?

- ▶ On peut introduire une marge d'erreur  $\xi_i$  pour la classification



# Hyperplan à marges douces

---

- ▶  $\xi_i = 0$  s'il n'existe pas d'erreur pour  $x_i$ 
  - ▶  $\xi_i$  sont des variables qui donnent du "mou" aux marges optimales

$$\begin{cases} \mathbf{w}^T \mathbf{x}_i + b \geq 1 - \xi_i & y_i = 1 \\ \mathbf{w}^T \mathbf{x}_i + b \leq -1 + \xi_i & y_i = -1 \\ \xi_i \geq 0 & \forall i \end{cases}$$

- ▶ Nous voulons minimiser

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

- ▶  $C$ : paramètre de compromis entre l'erreur et la marge
- ▶ Le problème d'optimisation devient

$$\begin{aligned} & \text{Minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ & \text{subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

# Détermination de l'hyperplan de séparation

---

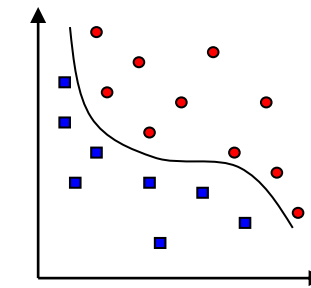
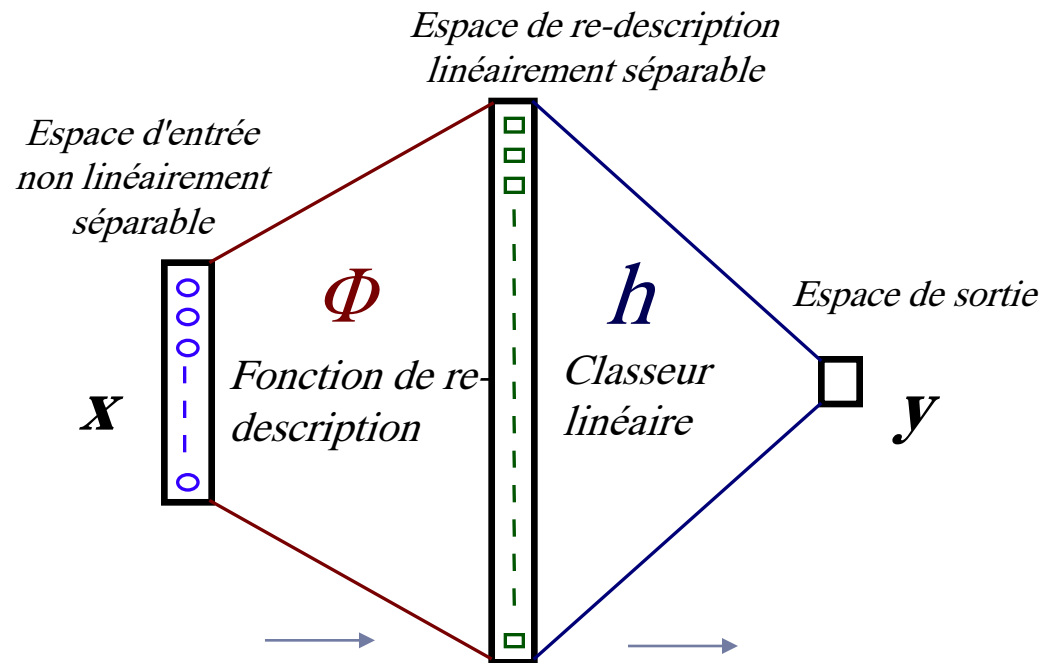
- ▶ La forme duale du problème est

$$\begin{aligned} \max. \quad W(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{subject to } C &\geq \alpha_i \geq 0, \quad \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

- ▶  $\mathbf{w}$  est aussi donné par  $\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}$
- ▶ La seule différence avec le cas linéairement séparable est qu'il existe une limite supérieure  $C$  aux  $\alpha_i$

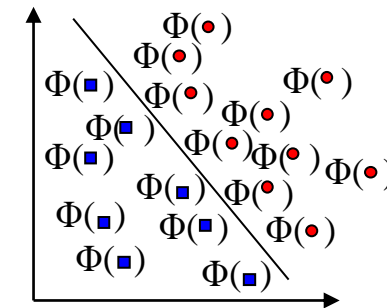
# Extension à une surface de séparation non-linéaire

- « Simplifier les choses » en projetant les  $x_i$  dans un espace qui les rend linéairement séparables



Espace d'entrée

$\Phi()$  ↓



Espace de re-description

# Modification due à la transformation

---

- ▶ Substituer les arguments transformés dans les produits scalaires lors de la phase d'apprentissage,

Problème original :

$$\max. W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

subject to  $C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$

Après transformation :

$$\max. W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$$

subject to  $C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$

- ▶ Cependant, trouver  $\Phi()$  n'est pas évident !

# Modification due à la transformation

---

- ▶ Les nouvelles données  $\mathbf{z}$  sont toujours classées dans la classe 1 si  $f \geq 0$ , la classe 2 sinon :

Original :

$$\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}$$
$$f = \mathbf{w}^T \mathbf{z} + b = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \mathbf{x}_{t_j}^T \mathbf{z} + b$$

Après transformation :

$$\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \phi(\mathbf{x}_{t_j})$$
$$f = \langle \mathbf{w}, \phi(\mathbf{z}) \rangle + b = \sum_{j=1}^s \alpha_{t_j} y_{t_j} \phi(\mathbf{x}_{t_j})^T \phi(\mathbf{z}) + b$$

et la surface de séparation dans le nouvel espace est :  $D(\mathbf{x}) = \sum_{j=1}^s \hat{\alpha}_j y_j \Phi(\mathbf{x}_j)^T \Phi(\mathbf{x}_j) + \hat{b}$

# Extension à une surface de séparation non-linéaire

---

- ▶ Problèmes cependant :

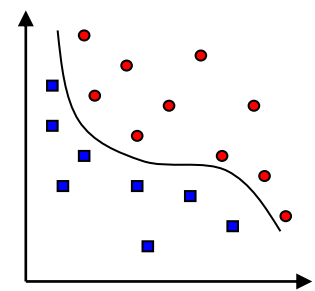
- ▶  $\Phi = ?$

- ▶ Grand effort de calcul potentiel ( $d$  explose !)

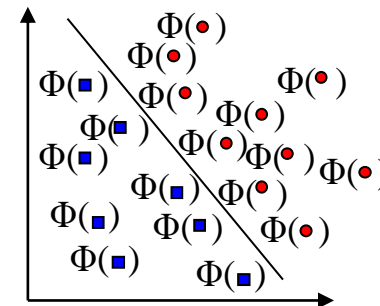
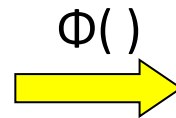
- ▶ SVM à fonctions noyaux résout les deux problèmes

- ▶ Efficacité computationnelle

- ▶ La transformation désirée des données est faite implicitement !



Espace d'entrée



Espace de re-description



# L'astuce des fonctions noyau

---

- ▶ Définition d'une fonction noyau :

$$K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$$

- ⇒ La connaissance de  $K(\cdot)$  permet de calculer un produit scalaire où intervient  $\Phi(\cdot)$  sans connaître l'expression de  $\Phi(\cdot)$
- ▶ Or, seuls des produits scalaires interviennent dans la solution du problème d'optimisation
  - ▶ Un autre avantage d'utiliser  $K(\cdot)$  est qu'il représente intuitivement la similarité entre les  $\mathbf{x}$  et  $\mathbf{y}$ , obtenue de nos connaissances a priori
  - ▶ Cependant,  $K(\mathbf{x}, \mathbf{y})$  doit satisfaire certaines conditions (conditions de Mercer) pour que le  $\Phi(\cdot)$  correspondant existe



# Les conditions de Mercer

---

- Pour une fonction  $K$  symétrique, il existe une fonction  $\Phi$  telle que :

$$K(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}') = \sum_{i=1}^m g_i(\mathbf{x}) \cdot g_i(\mathbf{x}')$$

ssi, pour toute fonction  $f$  telle que :  $\int f(\mathbf{x})^2 d\mathbf{x}$  est fini

l'on a :  $\int K(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \geq 0$

- Si cette condition est vérifiée, on peut appliquer la fonction noyaux dans le SVM
- **MAIS cela ne dit pas comment construire  $\Phi$**

# Exemple illustratif

---

- ▶ Définissons la fonction noyau  $K(\mathbf{x}, \mathbf{y})$  telle que, pour toute paire de vecteurs  $\mathbf{x}=(x_1, x_2)$  et  $\mathbf{y}=(y_1, y_2)$  :

$$K(\mathbf{x}, \mathbf{y}) = (1 + x_1 y_1 + x_2 y_2)^2$$

- ▶ Considérons maintenant une transformation  $\Phi$  qui prend un vecteur de dimension 2 et le projette dans un espace de dimension 6 :

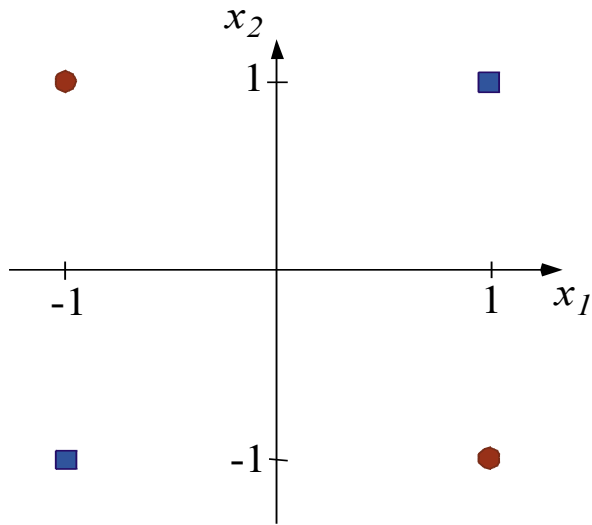
$$\Phi(\mathbf{x}) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

On peut voir en effectuant le calcul que

$$\begin{aligned} \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle &= (1 + x_1 y_1 + x_2 y_2)^2 \\ &= K(\mathbf{x}, \mathbf{y}) \end{aligned}$$

On peut donc obtenir le résultat sans avoir à passer par l'espace transformé!

# Illustration : le cas du XOR



Index $i$	$\mathbf{x}_i$	$y$
1	(1,1)	1
2	(1,-1)	-1
3	(-1,-1)	1
4	(-1,1)	-1

► Il faut résoudre :

$$\begin{cases} \max_{\alpha} \left( \sum_{i=1}^4 \alpha_i - \frac{1}{2} \sum_{i=1}^4 \sum_{j=1}^4 \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \right) \\ \forall i \quad 0 \leq \alpha_i \leq C \\ \sum_{i=1}^4 \alpha_i y_i = 0 \end{cases}$$

► Si on reprend la fonction noyau  $K(\mathbf{x}, \mathbf{y}) = (1 + x_1 y_1 + x_2 y_2)^2$ , on obtient les équations suivantes pour le Lagrangien :

$$\begin{aligned} Q(\alpha) = & \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 \\ & - \frac{1}{2} (9\alpha_1^2 - 2\alpha_1\alpha_2 + 2\alpha_1\alpha_3 - 2\alpha_1\alpha_4 \\ & + 9\alpha_2^2 - 2\alpha_2\alpha_3 + 2\alpha_2\alpha_4 + 9\alpha_3^2 - 2\alpha_3\alpha_4 + 9\alpha_4^2) \end{aligned}$$

$$\alpha_1 - \alpha_2 + \alpha_3 - \alpha_4 = 0$$

# Illustration : le cas du XOR

---

- ▶ Le maximum de  $Q(a)$  est obtenu en prenant ses dérivées par rapport aux  $\alpha_i$  et en trouvant les valeurs de  $\alpha_i$  qui les annulent :

$$\begin{cases} 1 - 9\alpha_1 + \alpha_2 - \alpha_3 + \alpha_4 = 0 \\ 1 + \alpha_1 - 9\alpha_2 + \alpha_3 - \alpha_4 = 0 \\ 1 - \alpha_1 + \alpha_2 - 9\alpha_3 - \alpha_4 = 0 \\ 1 + \alpha_1 - \alpha_2 + \alpha_3 - 9\alpha_4 = 0 \end{cases}$$

- La valeur optimale des multiplicateurs de Lagrange est :

$$\hat{\alpha}_1 = \hat{\alpha}_2 = \hat{\alpha}_3 = \hat{\alpha}_4 = \frac{1}{8}$$

- Les 4 données du ou exclusif sont donc des vecteurs de support, puisque aucune valeur trouvée de  $\alpha$  n'est nulle

# Illustration : le cas du XOR

- ▶ Dans l'espace de re-description :

$$\begin{cases} \hat{\mathbf{w}} = \sum_{i=1}^{n_s} \hat{\alpha}_i y_i \Phi(\mathbf{x}_i) \\ \hat{b} = y_s - \sum_{i=1}^{n_s} \hat{\alpha}_i y_i K(\mathbf{x}_i^T \mathbf{x}_s) \\ D(\mathbf{x}) = \sum_{j=1}^s \hat{\alpha}_j y_j K(\mathbf{x}_j, \mathbf{x}) + \hat{b} \end{cases}$$

- ▶ Donc :  $\hat{\mathbf{w}} = \frac{1}{8} [-\Phi(\mathbf{x}_1) + \Phi(\mathbf{x}_2) + \Phi(\mathbf{x}_3) - \Phi(\mathbf{x}_4)]$

$$= \frac{1}{8} \left[ - \begin{pmatrix} 1 \\ 1 \\ \sqrt{2} \\ 1 \\ -\sqrt{2} \\ -\sqrt{2} \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ -\sqrt{2} \\ 1 \\ -\sqrt{2} \\ \sqrt{2} \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ -\sqrt{2} \\ 1 \\ \sqrt{2} \\ -\sqrt{2} \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \\ \sqrt{2} \\ 1 \\ \sqrt{2} \\ \sqrt{2} \end{pmatrix} \right] = \begin{pmatrix} 0 \\ 0 \\ -1/\sqrt{2} \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

(on connaît  $\Phi()$  dans cet exemple, mais il n'est pas requis en général, car l'équation de la marge dépend seulement de  $K()$ )

$$\hat{b} = 1 - \frac{1}{8} \sum_{j=1}^4 y_j K(\mathbf{x}_j, \mathbf{x}_1) = 1 + \frac{1}{8} \sum_{j=1}^4 (-1)^j K(\mathbf{x}_j, \mathbf{x}_1) = 0$$

et :

$$D(\mathbf{x}) = \frac{1}{8} \sum_{j=1}^4 y_j K(\mathbf{x}_j, \mathbf{x}) = -\frac{1}{8} \sum_{j=1}^4 (-1)^j K(\mathbf{x}_j, \mathbf{x}) = -x_1 x_2$$

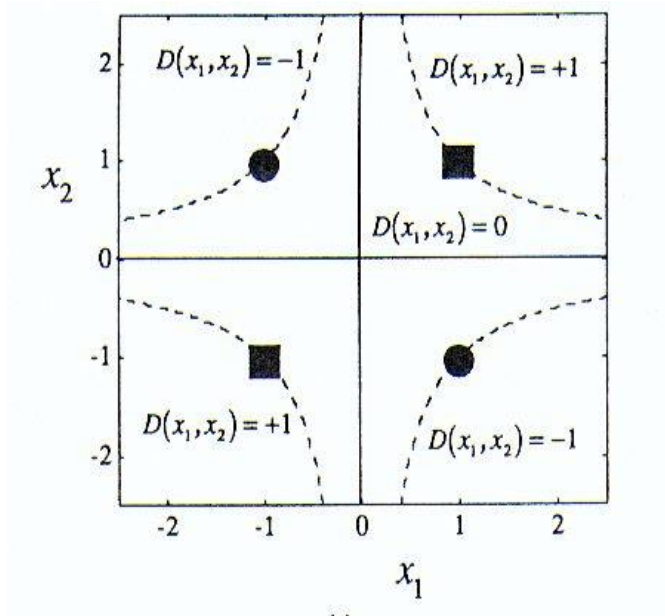
(on aurait obtenu le même résultat en utilisant  $\Phi()$  :

$$\hat{\mathbf{w}}^T \Phi(\mathbf{x}) = \left( 0, 0, \frac{-1}{\sqrt{2}}, 0, 0, 0 \right) \begin{pmatrix} 1 \\ x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \\ \sqrt{2} x_1 \\ \sqrt{2} x_2 \end{pmatrix} = -x_1 x_2$$

- ▶ La marge optimale est :

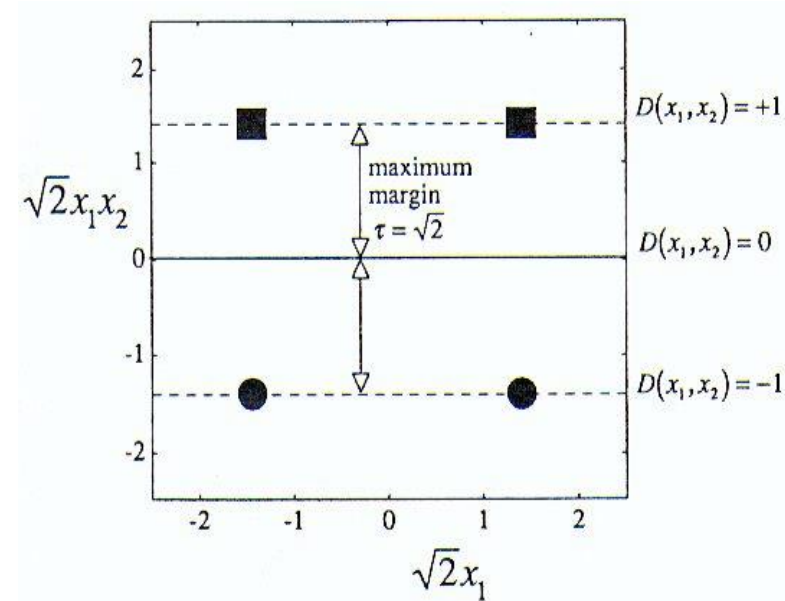
$$\begin{aligned} \frac{1}{2} \|\hat{\mathbf{w}}\|^2 &= \frac{1}{2} \hat{\mathbf{w}}^T \hat{\mathbf{w}} = \frac{1}{2} \left( \sum_{i=1}^4 \hat{\alpha}_i y_i \Phi(\mathbf{x}_i) \right)^T \left( \sum_{j=1}^4 \hat{\alpha}_j y_j \Phi(\mathbf{x}_j) \right) \\ &= \frac{1}{2} \cdot \frac{1}{8^2} \sum_{i=1}^4 \sum_{j=1}^4 y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) = \frac{1}{4} \Rightarrow \|\hat{\mathbf{w}}\| = \frac{1}{\sqrt{2}} \end{aligned}$$

# Illustration : le cas du XOR



Séparatrice dans l'espace d'entrée

$$D(x) = -x_1 x_2$$



Séparatrice dans l'espace  $\Phi(x)$

$$\sqrt{2} x_1 x_2 = 0$$

# Autre Exemple

---

- ▶ Supposons 5 nombres  $x_1=1, x_2=2, x_3=4, x_4=5, x_5=6$ , avec
  - ▶  $1, 2, 6 \in$  classe 1 ( $y=1$ )
  - ▶  $4, 5 \in$  classe 2 ( $y=-1$ )
  - ▶ Donc:  $\{(x_i, y_i)\}_{i=1,\dots,5} = \{(1,1), (2,1), (4,-1), (5,-1), (6,1)\}$
- ▶ Utilisons à nouveau le noyau polynomial de degré 2
  - ▶  $K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^\top \mathbf{y})^2$
  - ▶ C est choisi égal à 100
- ▶ Trouvons d'abord  $\{\alpha_i\}_{i=1,\dots,5}$  pour satisfaire :

$$\max. \sum_{i=1}^5 \alpha_i - \frac{1}{2} \sum_{i=1}^5 \sum_{j=1}^5 \alpha_i \alpha_j y_i y_j (x_i x_j + 1)^2$$
$$100 \geq \alpha_i \geq 0, \sum_{i=1}^5 \alpha_i y_i = 0$$

# Exemple

---

▶ La solution est :

▶  $\alpha_1=0, \alpha_2=2.5, \alpha_3=0, \alpha_4=7.333, \alpha_5=4.833$

▶ Les vecteur supports sont donc  $\{x_2=2, x_4=5, x_5=6\}$

▶ La fonction discriminante est

$$\begin{aligned} f(y) &= 2.5(1)(2y + 1)^2 + 7.333(-1)(5y + 1)^2 + 4.833(1)(6y + 1)^2 + b \\ &= 0.6667x^2 - 5.333x + b \end{aligned}$$

▶  $b$  trouvé en résolvant  $f(2)=1$  ou  $f(5)=-1$  ou  $f(6)=1$ , puisque  $x_2, x_4, x_5$  sont dans

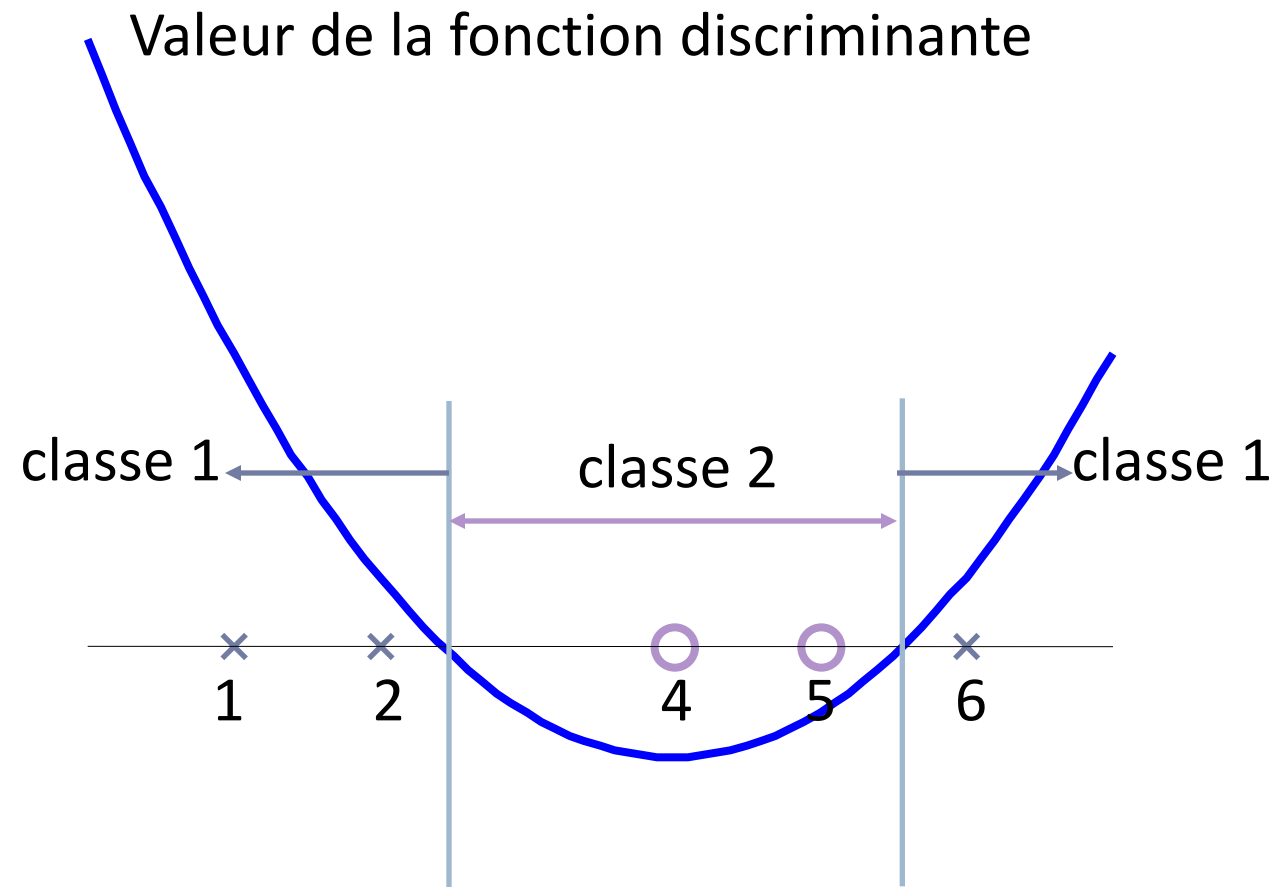
$y_i(\mathbf{w}^T \phi(z) + b) = 1$  et tous donnent  $b=9$

$$\longrightarrow f(y) = 0.6667x^2 - 5.333x + 9$$



# Exemple

---



# Exemples de fonctions noyaux

---

- ▶ Noyau polynomial de degré  $d$

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y} + 1)^d$$

- ▶ Noyau à fonction à base radiale de dispersion  $\sigma$

$$K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (2\sigma^2))$$

- ▶ Très proche des RN avec fonctions à base radiale
- ▶ Sigmoides avec paramètres  $\kappa$  et  $\theta$

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\kappa \mathbf{x}^T \mathbf{y} + \theta)$$

- ▶ Ne satisfait pas la condition de Mercer pour tous  $\kappa$  et  $\theta$
- ▶ La recherche d'autres fonctions noyaux pour diverses applications est très active !

# Exemple d'application avec SciKit

---

```
from sklearn import datasets
from sklearn.model_selection import train_test_split
from sklearn import svm
from sklearn import metrics

#Load dataset
cancer = datasets.load_breast_cancer()
# print the names of the 30 features
print("Features: ", cancer.feature_names)
# print the label type of cancer ('malignant' 'benign')
print("Labels: ", cancer.target_names)
# print data rows and features
print("rows, features per row: ", cancer.data.shape)

# Split dataset into 70% training set and 30% test set
X_train, X_test, y_train, y_test = train_test_split(cancer.data, cancer.target,
test_size=0.3,random_state=109)
```

```
#Create a svm Classifier
clf = svm.SVC(kernel='linear') # default=rbf
#Train the model using the training sets
clf.fit(X_train, y_train)
#Predict the response for test dataset
y_pred = clf.predict(X_test)

# Model Accuracy
print("Accuracy: ", metrics.accuracy_score(y_test, y_pred))
```

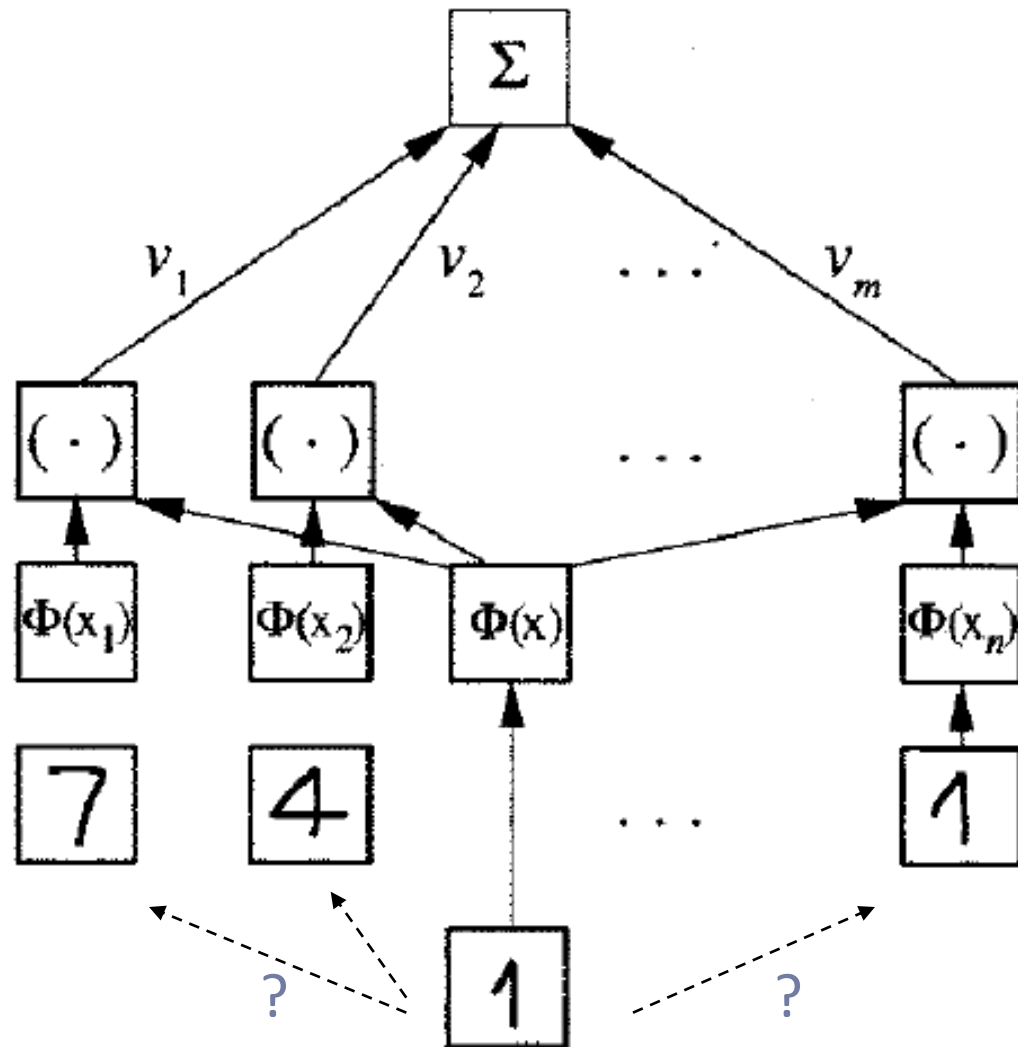


# Classification multi-classes

---

- ▶ SVM est à la base un classifieur binaire
- ▶ On peut changer la formulation pour permettre la classification multi-classe
  - ▶ L'ensemble des données est divisé en deux parts de multiples façons, et classé ensuite
    - ▶ Un contre tous ou un contre chaque alternative
    - ▶ Un SVM séparé est formé pour chaque division
  - ▶ La classification multi-classes est accomplie en combinant la sortie de tous les SVM
    - ▶ Prône à l'explosion combinatoire des possibilités!

# Exemple d'application des SVM : Reconnaissance de l'écriture manuscrite



output  $\Sigma v_i k(x, x_i) + b$

weights

dot product  $(\Phi(x) \cdot \Phi(x_i)) = k(x, x_i)$

mapped vectors  $\Phi(x_i), \Phi(x)$

support vectors  $x_1 \dots x_n$

test vector  $x$

## Sommaire: étapes de la classification

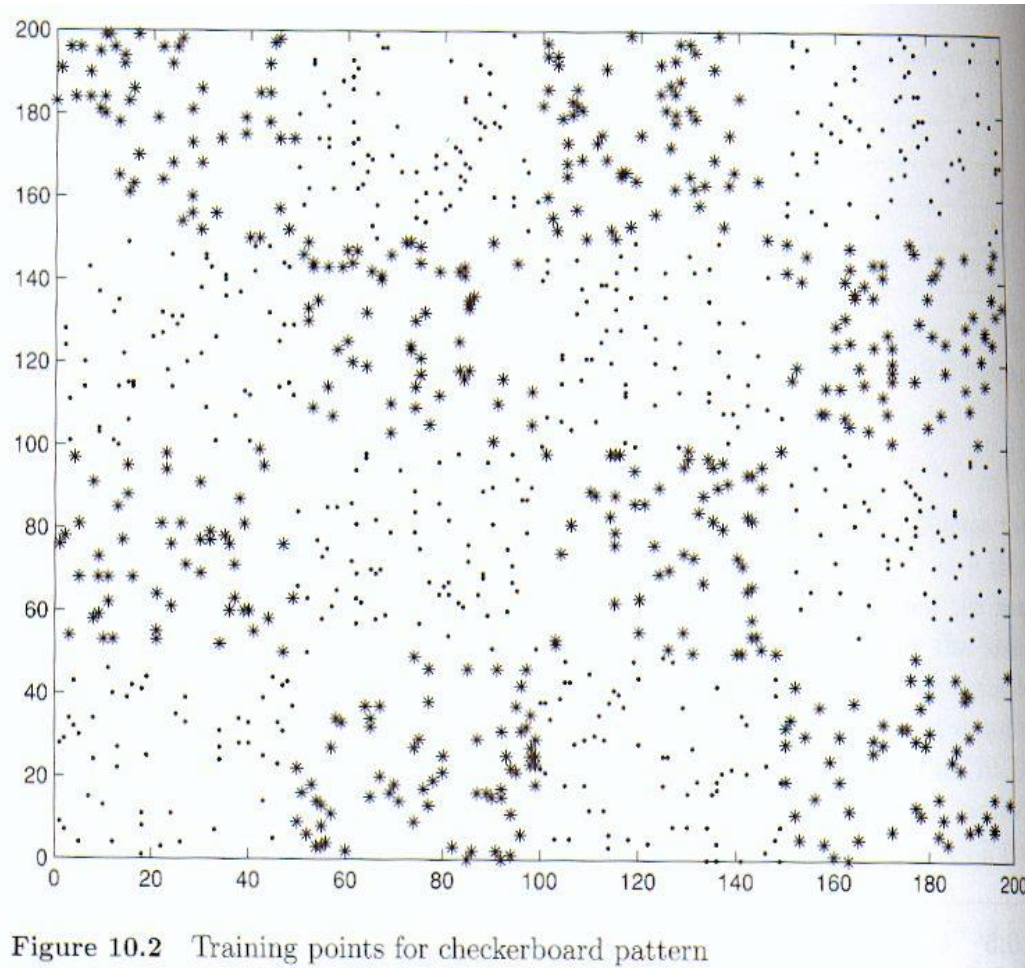
---

- ▶ Préparer la matrice des patrons
- ▶ Choisir la fonction noyau à utiliser
- ▶ Choisir les paramètres de la fonction noyau et la valeur de  $C$  (valeurs suggérées par le logiciel SVM ou essai-erreur).
- ▶ Exécuter l'algorithme d'apprentissage pour trouver  $\alpha_i$
- ▶ Les données nouvelles peuvent être classées en fonctions des  $\alpha_i$  et des vecteurs supports trouvés

# Effet des paramètres de contrôle.

- Apprentissage de données en damier
  - Apprentissage de deux classes
  - SVM à fonction noyau gaussienne

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

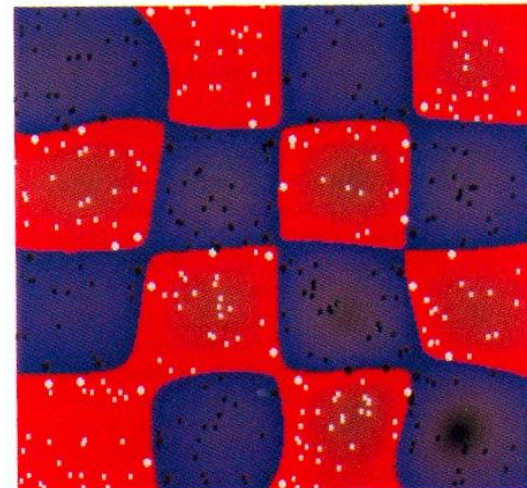
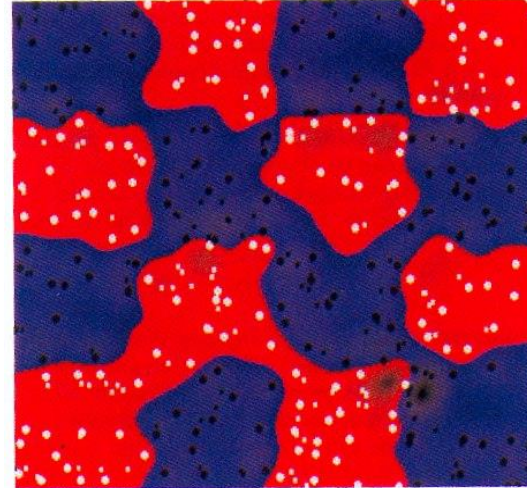


# Effet des paramètres de contrôle

- ▶ Apprentissage de deux classes
  - ▶ exemples tirés uniformément sur l'échiquier
- ▶ SVM à fonctions noyau gaussienne

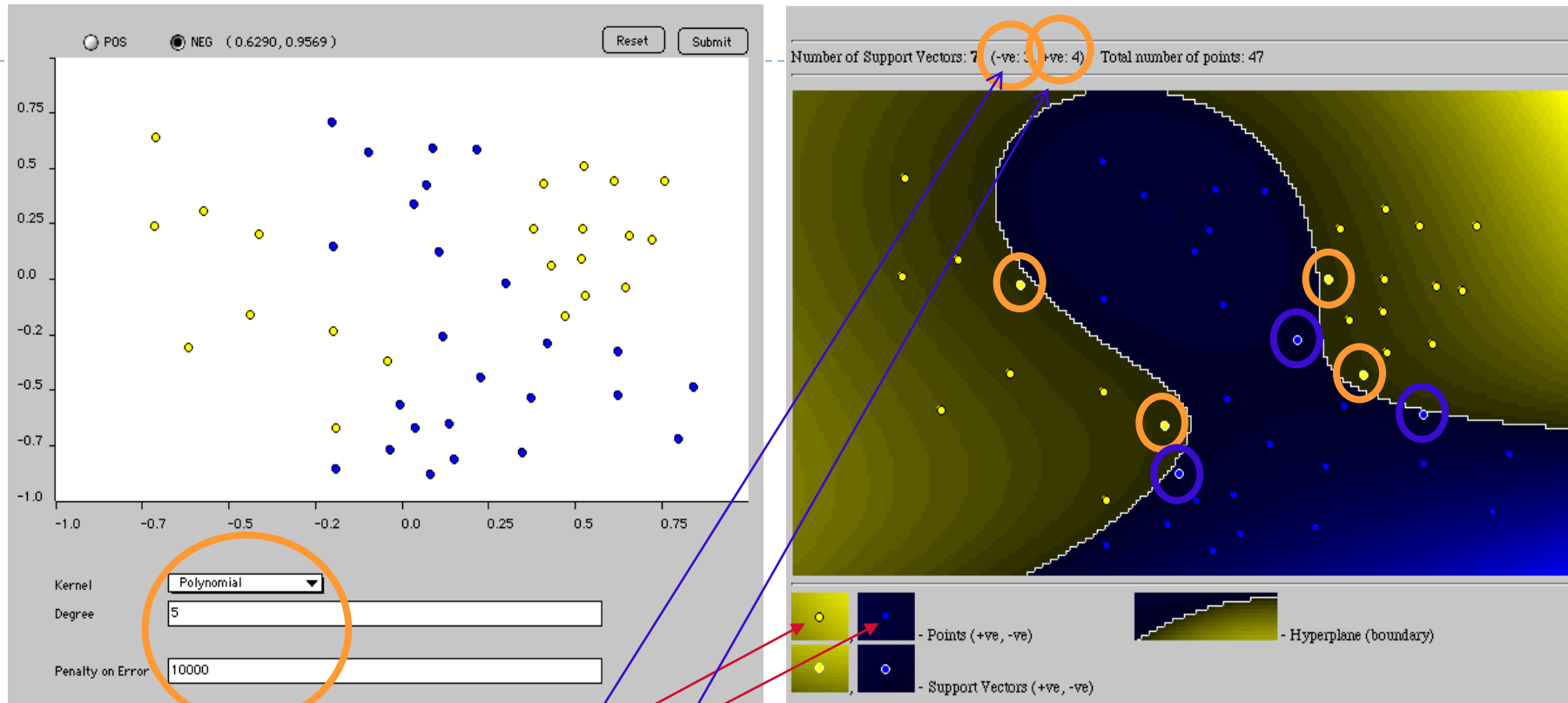
$$K(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}}$$

- ▶ Ici deux valeurs de  $\sigma$ 
  - ▶ En haut : petite valeur
  - ▶ En bas : grande valeur
- ▶ Les gros points sont des exemples critiques
  - ▶ Plus en haut qu'en bas



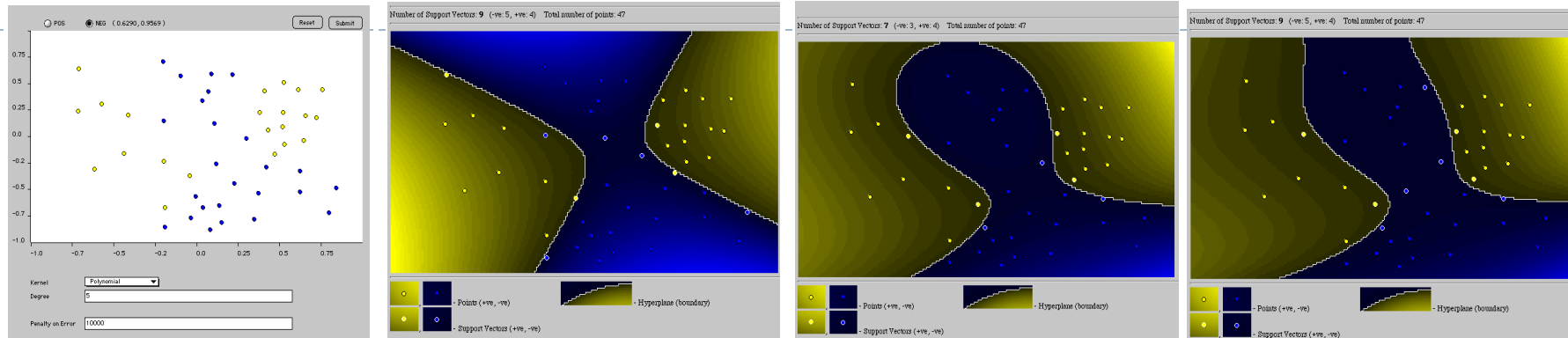


# Une applette de démonstration



- ▶ <http://sym.cs.rhul.ac.uk/pagesnew/GPat.shtml>
- ▶ 47 exemples (22 +, 25 -)
- ▶ *Exemples critiques*: 4 + et 3 -
- ▶ **Ici fonction polynomiale de degré 5 et  $C = 10000$**

# Paramètres de contrôle : les fonctions noyau



(5-, 4+)

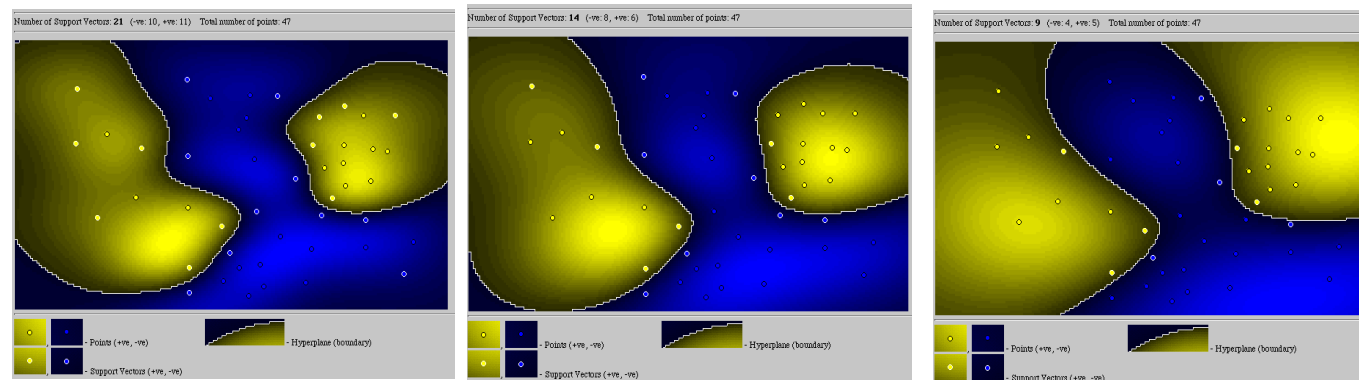
(3-, 4+)

(5-, 4+)

- ▶ 47 exemples (22 +, 25 -)

Ici *fonction polynomiale* de degré 2, 5, 8 et  $C = 10000$

- ▶ *Exemples critiques* : 4 + et 3 -



(10-, 11+)

(8-, 6+)

(4-, 5+)

Ici *fonction Gaussienne* de  $\sigma = 2, 5, 10, 20$  et  $C = 10000$

# Domaines d'application des SVMs

---

## ▶ Traitement d'images

- ▶ Reconnaissance de caractères manuscrits
- ▶ Reconnaissance de scènes naturelles
- ▶ Reconnaissance de visages

- ▶ *Entrées* : image bidimensionnelle en couleur ou en tons de gris codée en vecteur de pixels ou traits
- ▶ *Sortie* : classe (chiffre / personne)

# Application : images couleurs

---

- ▶ Ex. : Base d'images Corel Stock Photo Collection
  - ▶ 200 catégories
  - ▶ 100 images / catégorie
- ▶ Codage
  - ▶ Pixel = vecteur dans espace à trois dimensions (RGB)
  - ▶ Image = histogramme (fraction des pixels d'une couleur donnée)  
Invariant / nombreuses opérations

▶ Noyau :



$$K(x, z) = \exp\left(-\frac{d(x, z)}{\sigma^2}\right)$$

$$d(x, z) = \sum_{i=1}^n \frac{(x_i - z_i)^2}{x_i + z_i} \quad (\text{fonction } c^2)$$

# Domaines d'application des SVMs

---

## ▶ Catégorisation de textes

- ▶ Classification de courriels
- ▶ Classification de pages web
  
- ▶ *Entrées*: document (texte ou html)
  - ▶ Approche « sac de mots »
  - ▶ Document = vecteur de mots (lemmatisés pondérés par tf-idf)
- ▶ *Sortie*: catégorie (thème, spam/non-spam)
- ▶ ***Noyau***:
  - ▶ Produit scalaire des vecteurs
  - ▶  $C = \infty$  (marge dure)

# Domaines d'application des SVMs

---

## ▶ Diagnostic médical

- ▶ Évaluation du risque de cancer
  - ▶ Détection d'arythmie cardiaque
  - ▶ Évaluation du risque d'accidents cardio-vasculaires à moins de 6 ans
- 
- ▶ *Entrées* : état du patient (sexe, age, bilan sanguin, ...)
  - ▶ *Sortie* :
    - ▶ Classe : à risque ou non
    - ▶ Probabilité d'accident à échéance donnée

# Extensions

---

- ▶ Leçon à retenir des SVM:
  - ▶ *Un algorithme linéaire dans l'espace de re-description peut remplacer un algorithme non-linéaire dans l'espace d'entrée*
- ▶ Les algorithmes linéaires classiques peuvent être généralisés en des versions non-linéaires en allant vers l'espace de re-description
  - ▶ ACP à noyaux, k-moyennes à noyaux, etc.
  - ▶ Régression
  - ▶ Détection de « nouveautés »

# Régression vectorielle à support $\epsilon$ ( $\epsilon$ -SVR)

- ▶ L'idée est de trouver  $f(\vec{x}) = \vec{w} \cdot \vec{x} + b$  qui fait l'approximation d'un ensemble  $\{\vec{x}_i, y_i\}_{i=1..N}$  avec une déviation maximale de  $\pm\epsilon$  des vraies valeurs de  $y$ 
  - Les points doivent être compris dans la bande de largeur  $2\epsilon$ .

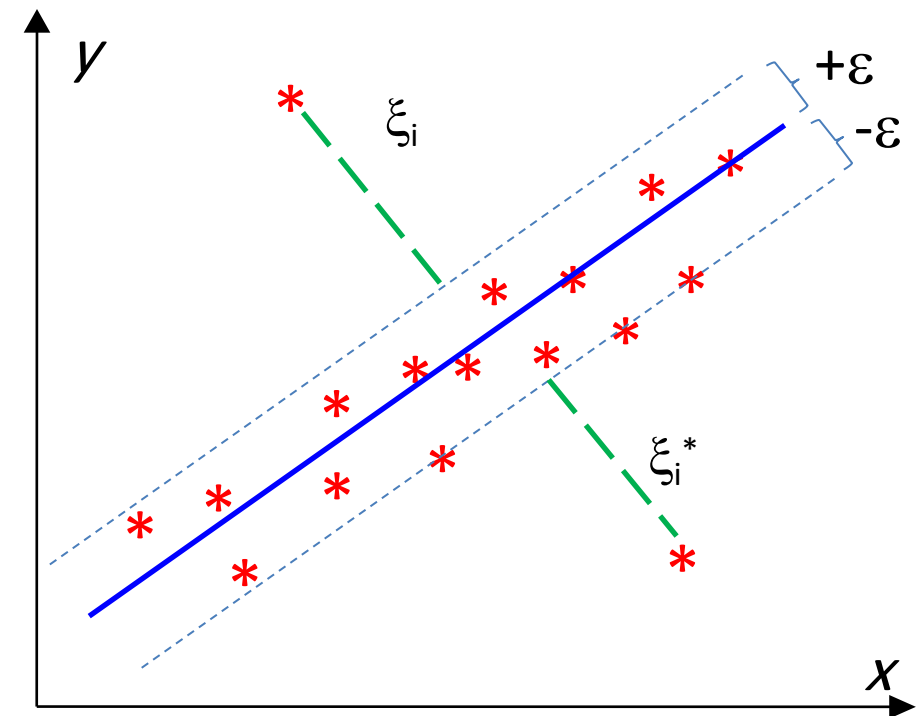
- ▶ Le problème d'optimisation est :

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

$$\text{subject to } \begin{cases} u_i - \mathbf{w}^T \mathbf{x}_i - b \leq \epsilon + \xi_i \\ \mathbf{w}^T \mathbf{x}_i + b - y_i \leq \epsilon + \xi_i^* \\ \xi_i \geq 0, \xi_i^* \geq 0 \end{cases}$$

ou  $\xi_i$  et  $\xi_i^*$  pénalisent les points hors cible

- ▶ Formulation similaire à SVM

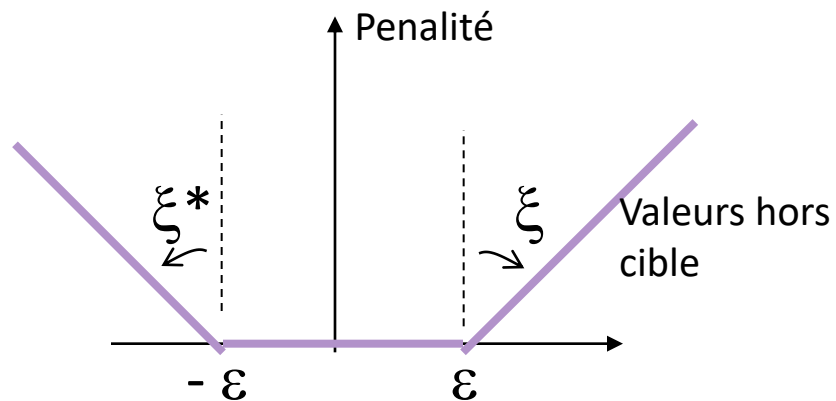




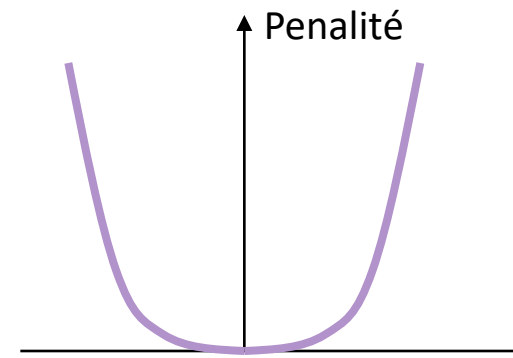
# Régression vectorielle à support $\varepsilon$ ( $\varepsilon$ -SVR)

- ▶ Régression linéaire dans l'espace de redescription
- ▶ À l'encontre de la régression par moindres carrés, la fonction d'erreur est une fonction de perte  $\varepsilon$ -insensible, et linéaire ensuite
  - ▶ Intuitivement, une erreur inférieure à  $\varepsilon$  est ignorée
  - ▶ Cela mène à des points de marge torses similaires à SVM

Fonction de perte  $\varepsilon$ -insensible



Fonction de perte quadratique



# Exemple d'application pour la regression

```
# Support Vector Regression (SVR) using linear and non-linear kernels
# 1D regression using linear, polynomial and RBF kernels.
import numpy as np
from sklearn.svm import SVR
import matplotlib.pyplot as plt

# Generate sample data
X = np.sort(5 * np.random.rand(40, 1), axis=0)
y = np.sin(X).ravel()

# Add noise to targets
y[::5] += 3 * (0.5 - np.random.rand(8))

# Fit regression model
svr_rbf = SVR(kernel='rbf', C=100, gamma=0.1, epsilon=.1)
svr_lin = SVR(kernel='linear', C=100, gamma='auto')
svr_poly = SVR(kernel='poly', C=100, gamma='auto', degree=3, epsilon=.1, coef0=1)

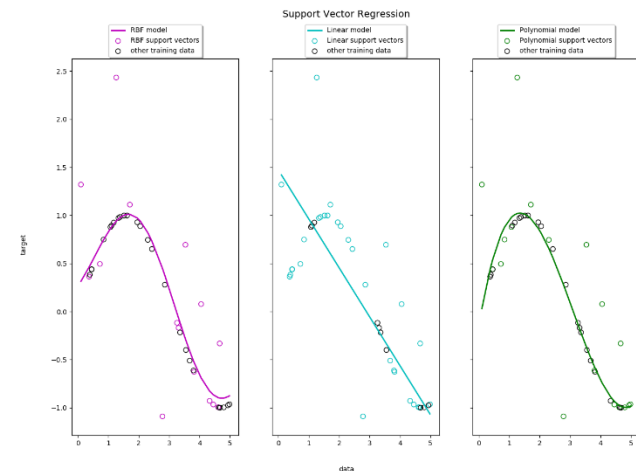
# Look at the results
lw = 2
svrs = [svr_rbf, svr_lin, svr_poly]
kernel_label = ['RBF', 'Linear', 'Polynomial']
model_color = ['m', 'c', 'g']

fig, axes = plt.subplots(nrows=1, ncols=3, figsize=(15, 10), sharey=True)
```

Source: scikit

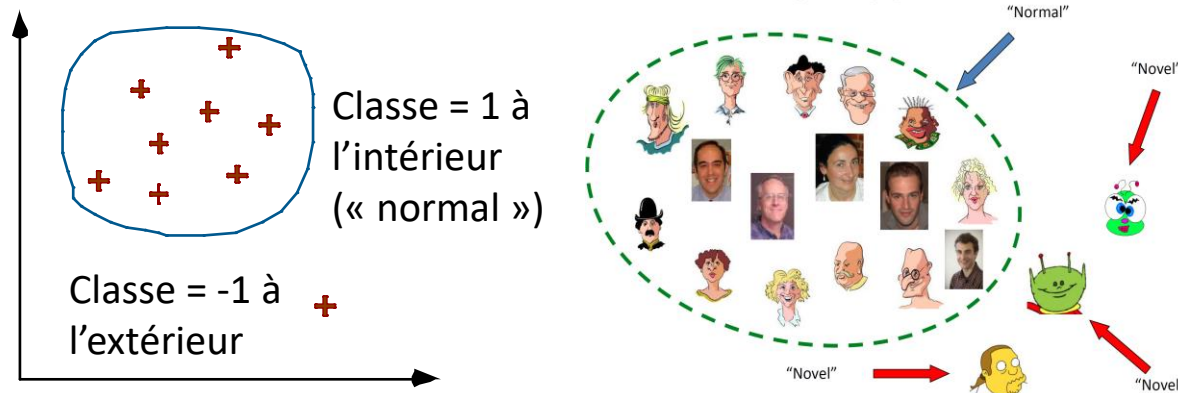
```
for ix, svr in enumerate(svrs):
    axes[ix].plot(X, svr.fit(X, y).predict(X), color=model_color[ix], lw=lw,
label='{}
                    model'.format(kernel_label[ix]))
    axes[ix].scatter(X[svr.support_], y[svr.support_], facecolor="none",
                    edgcolor=model_color[ix], s=50,
label='{} support
vectors'.format(kernel_label[ix]))
    axes[ix].scatter(X[np.setdiff1d(np.arange(len(X)), svr.support_)],
                    y[np.setdiff1d(np.arange(len(X)), svr.support_)], facecolor="none",
                    edgcolor="k", s=50, label='other training data')
    axes[ix].legend(loc='upper center', bbox_to_anchor=(0.5, 1.1),
                    fancybox=True, shadow=True)
ncol=1,
fig.text(0.5, 0.04, 'data', ha='center', va='center')
fig.text(0.06, 0.5, 'target', ha='left', va='middle')
fig.suptitle("Support Vec

plt.show()
```



# SVM pour le groupement à une classe (One Class SVM)

- ▶ Permet la détection de nouveautés ou d'anomalies
  - ▶ On cherche à séparer au maximum les points de la classe de l'origine



- ▶ On trouve  $f(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b)$  en minimisant :

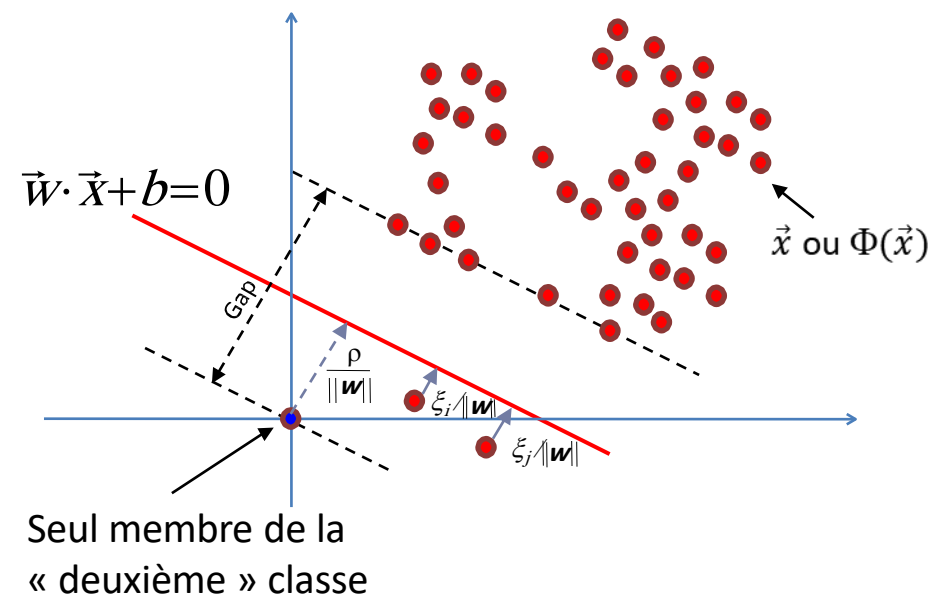
$$\frac{1}{2} \|\vec{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i + b$$

Sous les contraintes  $\vec{w} \cdot \vec{x} + b \geq -\xi_i$

$$\xi_i \geq 0 \text{ pour } i = 1, \dots, N$$

$\nu$  permet de régler la fraction de points singuliers

- ▶ Apprentissage non supervisé! (une seule classe)



# Exemple d'utilisation pour la détection d'anomalies

```
print(__doc__)
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.font_manager
from sklearn import svm

xx, yy = np.meshgrid(np.linspace(-5, 5, 500), np.linspace(-5, 5, 500))
# Generate train data
X = 0.3 * np.random.randn(100, 2)
X_train = np.r_[X + 2, X - 2]
# Generate some regular novel observations
X = 0.3 * np.random.randn(20, 2)
X_test = np.r_[X + 2, X - 2]
# Generate some abnormal novel observations
X_outliers = np.random.uniform(low=-4, high=4, size=(20, 2))
# fit the model
clf = svm.OneClassSVM(nu=0.1, kernel="rbf", gamma=0.1)
clf.fit(X_train)
y_pred_train = clf.predict(X_train)
y_pred_test = clf.predict(X_test)
y_pred_outliers = clf.predict(X_outliers)
n_error_train = y_pred_train[y_pred_train == -1].size
n_error_test = y_pred_test[y_pred_test == -1].size
n_error_outliers = y_pred_outliers[y_pred_outliers == 1].size
```

```
# plot the line, the points, and the nearest vectors to the plane
Z = clf.decision_function(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.title("Novelty Detection")
plt.contourf(xx, yy, Z, levels=np.linspace(Z.min(), 0, 7), cmap=plt.cm.PuBu)
a = plt.contour(xx, yy, Z, levels=[0], linewidths=2, colors='darkred')
plt.contourf(xx, yy, Z, levels=[0, Z.max()], colors='palevioletred')
s = 40
b1 = plt.scatter(X_train[:, 0], X_train[:, 1], c='white', s=s, edgecolors='k')
b2 = plt.scatter(X_test[:, 0], X_test[:, 1], c='blueviolet', s=s,
                edgecolors='k')
c = plt.scatter(X_outliers[:, 0], X_outliers[:, 1], c='gold', s=s,
                edgecolors='k')
plt.axis('tight')
plt.xlim((-5, 5))
plt.ylim((-5, 5))
plt.legend([a.collections[0], b1, b2, c],
          ["learned frontier", "training observations",
           "new regular observations", "new abnormal observations"],
          loc="upper left",
          prop=matplotlib.font_manager.FontProperties(size=11))
plt.xlabel(
    "error train: %d/200 ; errors novel regular: %d/40 ; "
    "errors novel abnormal: %d/40"
    % (n_error_train, n_error_test, n_error_outliers))
plt.show()
```

# SVM pour le groupement de plusieurs classes

- ▶ SVDD (support vector domain description) trouve l'hypersphère de plus petit rayon à englober les données

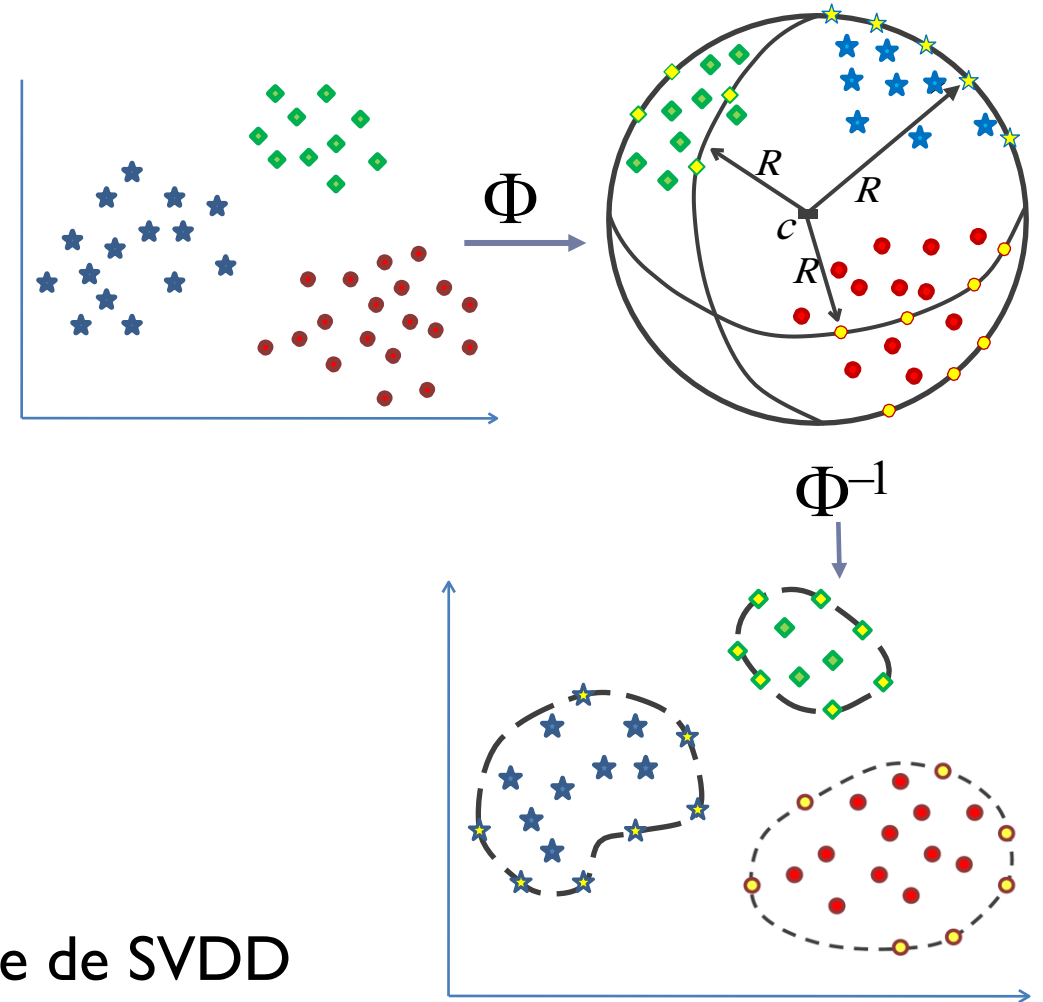
- ▶ Les vecteurs à la surface de l'hypersphère sont des vecteurs de support

- ▶ Pour  $X = \{x_i\}$  et  $\Phi$  une transformation de noyau radial (p. ex. gaussien), trouver la plus petite hypersphère de rayon  $R$  et centre  $c$  qui satisfait :

$$\|\Phi(x_j) - c\|^2 \leq R^2 + \xi_j$$

$$\xi_j \geq 0$$

- ▶ SVM à une classe avec noyau RBF est une forme de SVDD



# Pourquoi ça marche ?

---

## La marge est liée à la capacité en généralisation

▶ Normalement, la classe des hyperplans de  $\mathbb{R}^d$  est de  $d_H = d + 1$

▶ Mais la classe des hyperplans de marge  $\frac{1}{\|w\|}$  tq.  $\|w\|^2 \leq c$

est bornée par :  $d_H \leq \text{Min}(R^2 c, d) + 1$

où  $R$  est le rayon de la plus petite sphère englobant l'échantillon d'apprentissage  $S$

↙ Peut être beaucoup plus petit que la dimension  $d$  de l'espace d'entrée  $X$

# Forces et faiblesses des SVM

---

## ▶ Forces

- ▶ L'apprentissage est relativement facile
  - ▶ Pas de minima locaux comme pour les RNA
- ▶ L'algorithme est robuste face aux changements d'échelle
- ▶ Le compromis entre la complexité du classeur et l'erreur de classification peut être gérée explicitement
- ▶ Méthode générale
  - ▶ Des données non conventionnelles, telles des chaînes et des arbres peuvent servir d'entrées au SVM, à la place des vecteurs de traits
- ▶ Résultats en général **équivalents et souvent meilleurs**

## ▶ Faiblesses

- ▶ Il faut trouver la "bonne" fonction noyau
- ▶ Problèmes i.i.d. (données indépendantes et identiquement distribuées)
- ▶ Deux classes à la fois

# Sources documentaires

---

## ▶ Ouvrages / articles

- ▶ Cornuéjols & Miclet (02) : *Apprentissage artificiel. Concepts et algorithmes*. Eyrolles, 2002.
- ▶ Cristianini & Shawe-Taylor (00) : *Support Vector Machines and other kernel-based learning methods*. Cambridge University Press, 2000.
- ▶ Herbrich (02) : *Learning kernel classifiers*. MIT Press, 2002.
- ▶ Schölkopf, Burges & Smola (eds) (98) : *Advances in Kernel Methods : Support Vector Learning*. MIT Press, 1998.
- ▶ Schölkopf & Smola (02) : *Learning with kernels*. MIT Press, 2002.
- ▶ Smola, Bartlett, Schölkopf & Schuurmans (00) : *Advances in large margin classifiers*. MIT Press, 2000.
- ▶ Vapnik (95) : *The nature of statistical learning*. Springer-Verlag, 1995.

## ▶ Sites web

- ▶ <http://www.kernel-machines.org/> **(point d'entrée)**
- ▶ <http://www.support-vector.net> **(point d'entrée)**



# Implémentation des SVMs

---

- ▶ **Minimisation de fonctions différentiables convexes à plusieurs variables**
  - ▶ Pas d'optima locaux
  - ▶ **Mais :**
    - ▶ **Problèmes de stockage de la matrice noyau (si milliers d'exemples)**
    - ▶ **Long dans ce cas**
  - ▶ D'où mise au point de méthodes spécifiques
    - ▶ Gradient sophistiqué
    - ▶ Méthodes itératives, optimisation par morceaux
  - ▶ **Plusieurs packages publics disponibles**
    - ▶ SVMTorch
    - ▶ SVM<sup>Light</sup>
    - ▶ SMO
    - ▶ ...

# Logiciels

---

- ▶ Une liste de réalisations de SVM se trouve à <http://www.kernel-machines.org/software.html>
- ▶ Certaines (tel LIBSVM) peuvent gérer la classification multi-classe
- ▶ SVMLight figure parmi les premières mises en oeuvres de SVM (écrit en c ; <http://svmlight.joachims.org/>)
- ▶ IL existe plusieurs boîtes à outils Matlab pour les SVM sur le web
- ▶ Scikit aussi inclut SVM

# Autres ressources

---

- ▶ <http://www.kernel-machines.org/>
- ▶ <http://www.support-vector.net/>
- ▶ <http://www.support-vector.net/icml-tutorial.pdf>
- ▶ <http://www.kernel-machines.org/papers/tutorial-nips.ps.gz>
- ▶ <http://www.clopinet.com/isabelle/Projects/SVM/applist.html>

# Conclusion

---

- ▶ SVM sont une alternative aux réseaux de neurones pour la classification binaire
- ▶ Concepts clés des SVM : maximiser la marge et exploiter l'astuce des noyaux
- ▶ Domaine de recherche encore d'intérêt
- ▶ Plusieurs mises en œuvre en logiciel libre existent et sont disponibles sur le WEB !