

Apprentissage, réseaux de neurones et modèles graphiques

Chapitre 7 : Machines à vecteurs support

Michel Crucianu

avec des contributions de Jean-Philippe Tarel

<http://cedric.cnam.fr/~crucianm/ml.html>

25 septembre 2013

RCP209

1

Contenu du chapitre

- SVM pour la discrimination
 - ◆ Maximisation de la marge
 - ◆ Astuce des noyaux (*kernel trick*)
- Estimation du support d'une distribution
- SVM pour la régression
- Méthodes à noyaux et analyse des données
- Ingénierie des noyaux
 - ◆ Noyaux valides et « bons » noyaux
 - ◆ Noyaux pour différents types de données
 - ◆ Combinaison de noyaux, noyaux hybrides
- Exemples d'utilisation des SVM

25 septembre 2013

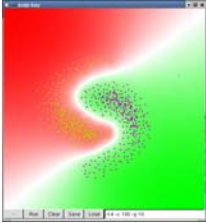
RCP209

2

le cnam

SVM pour la discrimination

- *Support Vector Machines*
 1. Séparateurs à Vastes Marges
 - ◆ SVM ne servent pas seulement à séparer (aussi à la régression, ...)
 - ◆ SVM pas les seuls séparateurs à vastes marges (autre ex. : *boosting*)
 2. Machines à vecteurs support



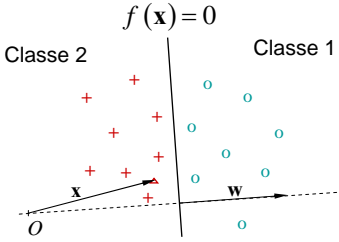
Exemple simple avec noyau angulaire
 Intensité de la couleur proportionnelle à l'éloignement de la frontière
 (outil employé : version maison de svm-toy)

25 septembre 2013
RCP209
3

le cnam

Classes linéairement séparables

- Données d'apprentissage $D_n = \{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq n}$, $y_i \in \{-1, +1\}$
- On cherche une fonction de décision $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$, basée sur l'hyperplan $\mathbf{w}^T \mathbf{x} + b = 0$, \mathbf{w} étant un vecteur normal à l'hyperplan
 - ◆ Affectation à la classe 1 si $f(\mathbf{x}) > 0$
 - ◆ Affectation à la classe 2 si $f(\mathbf{x}) < 0$
- **Séparabilité linéaire** : $\exists \mathbf{w}, b$ tels que $y_i f(\mathbf{x}_i) > 0$ pour $1 \leq i \leq n$
- Remarque : si la condition est valable pour \mathbf{w}, b , alors elle est valable pour $k\mathbf{w}, kb, \forall k > 0$

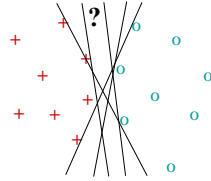


25 septembre 2013
RCP209
4

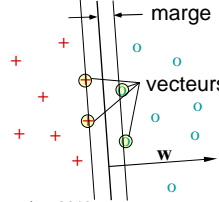
le cnam

Discrimination linéaire et marge

- Quelle séparation choisir quand plusieurs (éventuellement une infinité) sont possibles ?



- Une possibilité : choisir la séparation qui **maximise la marge** (distance minimale entre un exemple et la surface de séparation)



en fonction de \mathbf{w} : $\text{marge} = |\mathbf{x}_s^T \mathbf{w} + b| / \|\mathbf{w}\|$

normalisation : $|\mathbf{x}_s^T \mathbf{w} + b| = 1$ pour tout **vecteur support** \mathbf{x}_s

donc $\text{marge} = 1 / \|\mathbf{w}\|$

25 septembre 2013
RCP209
5

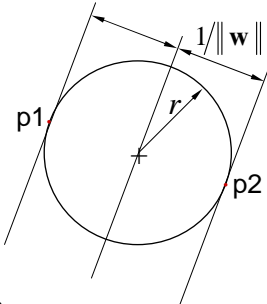
le cnam

Pourquoi s'intéresser à la marge

1. Marge et VC-dimension [Vap98] : pour des données qui se situent à l'intérieur d'une (hyper)sphère de rayon r , la VC-dimension h de la famille de (hyper)plans de marge $1/\|\mathbf{w}\| \geq m$ est bornée,

$$h \leq (r/m)^2 + 1$$

Exemple : dans \mathbb{R}^2 la VC-dimension des droites est $h = 3$, mais la VC-dimension des droites de marge $1/\|\mathbf{w}\| \geq r$ est $h = 2$ pour des données du disque de rayon r

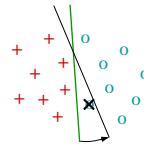


25 septembre 2013
RCP209
6

Pourquoi s'intéresser à la marge (2)

2. Vecteurs support et *leave one out cross-validation* [SS02]

- Validation croisée *leave one out* : on retire un exemple de D_n et on apprend sur les $n-1$ restants, on évalue l'erreur sur l'exemple retiré, on répète la procédure pour chaque exemple de D_n ; la moyenne des erreurs ainsi obtenues constitue une estimation du risque espéré
- Le risque espéré a comme borne supérieure $E[S_{(n)}] / n$, où $E[S_{(n)}]$ est l'espérance du nombre de vecteurs support obtenus en maximisant la marge sur tous les choix possibles d'ensembles d'apprentissage de taille n



Problème d'optimisation

- Maximiser la marge : optimisation sous contraintes d'inégalité (problème primal)

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad 1 \leq i \leq n \end{cases}$$

- En introduisant les multiplicateurs de Lagrange α

$$L(\mathbf{w}, b, \alpha) = -\frac{1}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^n \alpha_i [y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1]$$

- Des conditions nécessaires :

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

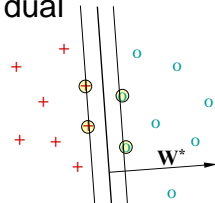
$$\frac{\partial L}{\partial \mathbf{w}}(\mathbf{w}^*, b^*, \alpha^*) = 0 \Rightarrow \mathbf{w}^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$$

le cnam

Résolution du problème d'optimisation

- Par substitution on obtient le problème dual

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \alpha_i \geq 0, \quad 1 \leq i \leq n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$



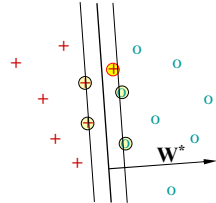
- Les vecteurs support sont ceux pour lesquels $\alpha_i \neq 0$
- b^* est obtenu à partir de la condition $|\mathbf{x}_s^T \mathbf{w}^* + b^*| = 1$ valable pour les vecteurs support
- La fonction de décision permettant de classer une nouvelle observation \mathbf{x} sera $f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i^T \mathbf{x} + b^*$

25 septembre 2013 RCP209 9

le cnam

Cas des données non séparables

- Si les données d'apprentissage ne sont pas linéairement séparables, on emploie les contraintes assouplies $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$ pour $1 \leq i \leq n$ et avec $\xi_i \geq 0$
- Le nouveau problème d'optimisation doit inclure une pénalité pour les erreurs ($\xi_i > 0$) faites sur l'ensemble d'apprentissage

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad 1 \leq i \leq n, \quad \xi_i \geq 0 \end{cases}$$


25 septembre 2013 RCP209 10

Cas des données non séparables (2)

- Le problème dual résultant est

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ C \geq \alpha_i \geq 0, \quad 1 \leq i \leq n \\ \sum_{i=1}^n \alpha_i y_i = 0 \end{cases}$$

- C est une constante de régularisation (la régularisation est d'autant plus forte que C est proche de 0 !)
- Pour une meilleure stabilité numérique, b^* est obtenu à partir de la moyenne sur l'ensemble I des vecteurs pour lesquels $0 < \alpha_i < C$

$$b^* = (1/|I|) \sum_{i \in I} (y_i - \sum_{j=1}^n \alpha_j^* y_j \mathbf{x}_i^T \mathbf{x}_j)$$

- La fonction de décision permettant de classer une nouvelle observation \mathbf{x} est toujours

$$f^*(\mathbf{x}) = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i^T \mathbf{x} + b^*$$

25 septembre 2013

RCP209

11

Conditions KKT

- Conditions nécessaires et parfois suffisantes d'optimalité (Karush-Kuhn-Tucker) pour le problème dual :

$$\alpha_i = 0 \Rightarrow y_i f(\mathbf{x}_i) > 1 \text{ et } \xi_i = 0 \leftarrow \text{hors marge}$$

$$0 < \alpha_i < C \Rightarrow y_i f(\mathbf{x}_i) = 1 \text{ et } \xi_i = 0 \leftarrow \text{sur la marge}$$

$$\alpha_i = C \Rightarrow y_i f(\mathbf{x}_i) < 1 \text{ et } \xi_i \geq 0 \leftarrow \text{dans la marge ou mal classés}$$

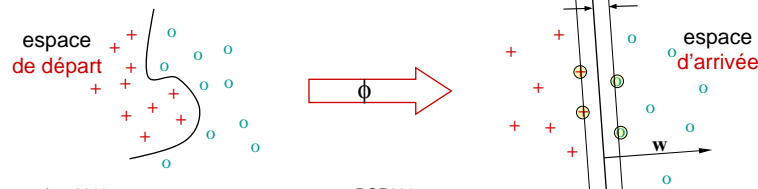
25 septembre 2013

RCP209

12

Astuce des noyaux (*kernel trick*)

- Comment étendre ces résultats à des séparateurs non linéaires ?
- Principe : transposer les données dans un autre espace (en général de plus grande dimension) dans lequel elles sont linéairement séparables (ou presque)
- Transformation $\phi: \mathbb{R}^d \rightarrow \mathcal{H}$, $\mathbf{x} \rightarrow \phi(\mathbf{x})$, \mathcal{H} espace de Hilbert
- Sous certaines conditions, l'existence de ϕ et de \mathcal{H} est garantie et ϕ est associée à une fonction noyau $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$, telle que $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j) \Rightarrow$ calculs faits dans l'espace de départ !



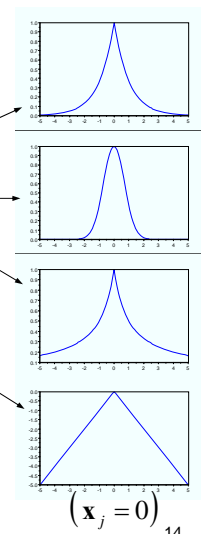
25 septembre 2013

RCP209

13

Exemples de noyaux

- Linéaire : $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Exponentiel : $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|)$
- Gaussien (RBF) : $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2)$
- Hyperbolique : $K(\mathbf{x}_i, \mathbf{x}_j) = 1 / (\varepsilon + \gamma \|\mathbf{x}_i - \mathbf{x}_j\|)$
- Angulaire : $K(\mathbf{x}_i, \mathbf{x}_j) = -\|\mathbf{x}_i - \mathbf{x}_j\|$
- Puissance : $K(\mathbf{x}_i, \mathbf{x}_j) = -\|\mathbf{x}_i - \mathbf{x}_j\|^\beta$ ($0 < \beta \leq 2$)



25 septembre 2013

RCP209

14

Exemple de noyau polynomial

- Noyau polynomial de degré 2 : $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$
- Dans un espace 2D, $\mathbf{x}_i = (x_i, y_i)$
- Donc, en développant,

$$K(\mathbf{x}_i, \mathbf{x}_j) = \begin{pmatrix} 1 \\ \sqrt{2}x_i \\ \sqrt{2}y_i \\ x_i^2 \\ \sqrt{2}x_i y_i \\ y_i^2 \end{pmatrix}^T \begin{pmatrix} 1 \\ \sqrt{2}x_j \\ \sqrt{2}y_j \\ x_j^2 \\ \sqrt{2}x_j y_j \\ y_j^2 \end{pmatrix}$$

→ Produit scalaire en dimension 6 !

Résolution du problème d'optimisation

- Dans la formulation du problème d'optimisation, les produits scalaires $\mathbf{x}_i^T \mathbf{x}_j$ sont remplacés par $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j)$, donc pour des données qui ne sont pas complètement séparables dans l'espace \mathcal{H} le problème primal est

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ y_i (\langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b) \geq 1 - \xi_i, 1 \leq i \leq n, \xi_i \geq 0 \end{cases}$$

et le problème dual correspondant :

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

avec b^* donné par

$$b^* = \frac{1}{|I|} \sum_{i \in I} \left(y_i - \sum_{j=1}^n \alpha_j^* y_j K(\mathbf{x}_i, \mathbf{x}_j) \right)$$

et la fonction de décision

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

le cnam

Effet de l'échelle pour noyau RBF

$1/\gamma=0.001$, 27 SV

$1/\gamma=0.01$, 26 SV

$1/\gamma=0.1$, 12 SV

$1/\gamma=1.0$, 10 SV

Sur-apprentissage :
mauvaise généralisation,
apprentissage du bruit

Sous-apprentissage :
mauvaise généralisation,
frontière imprécise

25 septembre 2013

RCP209

17

le cnam

Effet de l'échelle et de C pour RBF

Noyau angulaire : C=100 (gauche) et 10 (droite) Noyau RBF $\gamma=100$: C=100 (gauche) et 10 (droite)

Noyau RBF avec $\gamma=20$ & C=100, puis $\gamma=100$ & C=100, puis $\gamma=1000$ & C=100, puis $\gamma=1000$ & C=0,7

25 septembre 2013

RCP209

18

Formulation alternative : ν -SVM

- La signification de la constante C n'est pas très intuitive, on reformule le problème d'optimisation de la façon suivante

$$\begin{cases} \min_{\alpha} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ 1/n \geq \alpha_i \geq 0, \quad 1 \leq i \leq n \\ \sum_{i=1}^n \alpha_i y_i = 0 \\ \sum_{i=1}^n \alpha_i \geq \nu \end{cases}$$

- Sous certaines conditions, $\nu \in (0, 1]$ est
 - ◆ Une borne supérieure pour la fraction d' « erreurs de marge »
 - ◆ Une borne inférieure pour la fraction de vecteurs de support


Types de bases

- Base d'**apprentissage** (*training database*) : données étiquetées permettant de construire le modèle de l'objet d'intérêt, échantillonnées (en général) pour être représentatives
- Base de **validation** (*validation database*) : données étiquetées différentes de la base d'apprentissage, permettant d'optimiser le choix des paramètres (C, noyau, etc.)
- Base de **test** (*test database*) : données étiquetées différentes des bases précédentes, permettant d'évaluer les performances du modèle appris

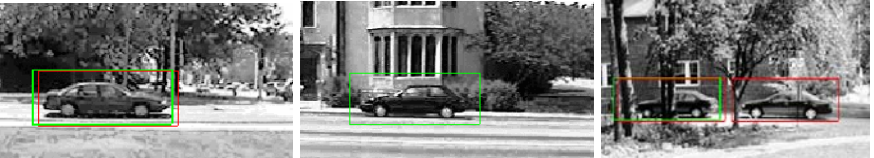
le cnam

Exemple applicatif : détection véhicule

- Base d'apprentissage



- Résultats sur base de test



25 septembre 2013
RCP209
21

le cnam

Courbes ROC : sensibilité, spécificité

- Évaluation de la discrimination entre présence d'une classe et absence de la classe (ex. : détection d'objets)
- Table de vérité

		Vérité terrain	
		Classe présente	Classe absente
Détecteur	Classe détectée	Vrai Positif	Faux Positif
	Classe non détectée	Faux Négatif	Vrai Négatif

- Sensibilité et spécificité pour un seuil de décision fixé

$$TPR = \text{Sensibilité} = \frac{VP}{\text{Total positif}} = \frac{VP}{VP + FN}$$

$$FPR = 1 - \text{Spécificité} = \frac{FP}{\text{Total négatif}} = \frac{FP}{VN + FP} = 1 - \frac{VN}{VN + FP}$$

25 septembre 2013
RCP209
22

le cnam

Évaluation comparative

- Courbe ROC (*Receiver Operating Characteristics*) : sensibilité vs. 1-spécificité pour différents seuils de décision
- En cas d'intersection des courbes ROC choisir un point sur la courbe ROC par exemple avec l'indice Dice maximum

$$\text{Dice} = \frac{2 VP}{2 VP + FP + FN}$$

25 septembre 2013
RCP209
23

le cnam

Support d'une distribution

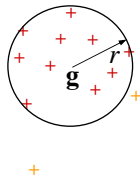
- Soit $\mathcal{D} = \{x_1, \dots, x_n\}$ un ensemble d'observations dans \mathcal{X} , issues de variables i.i.d. suivant la densité de probabilité $p(x)$ (inconnue)
- On cherche à estimer le **support** de cette densité ← moins de difficultés que pour l'estimation de la densité

Exemple simple avec noyau RBF
 Intensité de la couleur proportionnelle à l'éloignement de la frontière
 (outil employé : version maison de svm-toy)

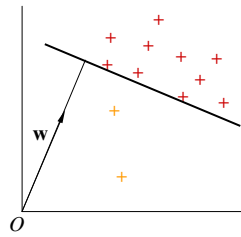
25 septembre 2013
RCP209
24

One-class SVM : alternatives

1. Trouver dans l'espace d'arrivée \mathcal{H} la plus petite hyper-sphère englobant les données
2. Trouver dans l'espace d'arrivée \mathcal{H} l'hyperplan le plus éloigné de l'origine, qui sépare les données de l'origine



Formulation 1 [TD99] : hyper-sphère englobant les données



Formulation 2 [SPS99] : hyperplan séparant de l'origine

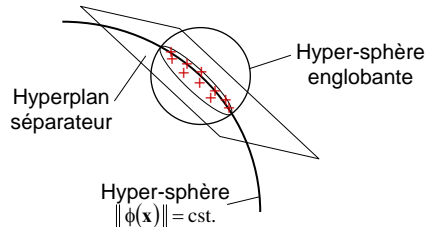
25 septembre 2013

RCP209

25

Condition d'équivalence

- Les deux formulations sont équivalentes si le noyau employé satisfait la condition $K(x, x) = R^2(\text{const.}), \forall x \in \mathcal{X}$
 - ◆ On a alors $\|\phi(x)\| = R$, donc l'image de l'espace de départ dans l'espace d'arrivée se situe sur l'hyper-sphère de rayon R
 - ◆ L'hyper-sphère de rayon r minimal ou l'hyperplan le plus éloigné de l'origine sépareront le même segment de l'hyper-sphère de rayon R :



25 septembre 2013

RCP209

26

One-class SVM : formulation 1

- En notant par \mathbf{g} le centre et par r le rayon de l'hyper-sphère et en faisant appel aux « variables d'assouplissement » $\xi_i \geq 0$, le problème d'optimisation à résoudre est

$$\begin{cases} \min_{\mathbf{g}, r, \xi} \left(r^2 + (1/\nu n) \sum_{i=1}^n \xi_i \right) \\ \|\phi(\mathbf{x}_i) - \mathbf{g}\|^2 \leq r^2 + \xi_i, 1 \leq i \leq n, \xi_i \geq 0 \end{cases}$$

- Avec les multiplicateurs de Lagrange on obtient le problème dual

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}_i) - \sum_{i,j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) & \text{avec } \mathbf{g} = \sum_{i=1}^n \alpha_i \phi(\mathbf{x}_i) \\ (1/\nu n) \geq \alpha_i \geq 0, 1 \leq i \leq n \\ \sum_{i=1}^n \alpha_i = 1 \end{cases}$$

Un nouveau point \mathbf{x} appartient au support si $\|\phi(\mathbf{x}) - \mathbf{g}\|^2 \leq r^2$, ou

$$K(\mathbf{x}, \mathbf{x}) - 2 \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i) - \sum_{i,j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \leq r^2$$

25 septembre 2013

RCP209

27

One-class SVM : formulation 2

- Avec la fonction de décision $f(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle - \rho)$ en faisant appel aux « variables d'assouplissement » $\xi_i \geq 0$, avec $\nu \in (0, 1]$, le problème d'optimisation à résoudre est

$$\begin{cases} \min_{\mathbf{w}, \xi, \rho} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu n} \sum_{i=1}^n \xi_i - \rho \\ \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle \geq \rho - \xi_i, 1 \leq i \leq n, \xi_i \geq 0 \end{cases}$$

- Avec les multiplicateurs de Lagrange on obtient le problème dual

$$\begin{cases} \min_{\alpha} \sum_{i,j=1}^n \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) & \text{avec } \rho = \langle \mathbf{w}, \phi(\mathbf{x}_s) \rangle = \sum_{i=1}^n \alpha_i K(\mathbf{x}_i, \mathbf{x}_s) \\ (1/\nu n) \geq \alpha_i \geq 0, 1 \leq i \leq n & \text{pour tout } \mathbf{x}_s \text{ tel que } (1/\nu n) > \alpha_i > 0 \\ \sum_{i=1}^n \alpha_i = 1 & \text{Un nouveau point } \mathbf{x} \text{ appartient} \\ & \text{au support si } \text{sgn}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle - \rho) = 1 \end{cases}$$

25 septembre 2013

RCP209

28

One-class SVM : formulation 2 (2)

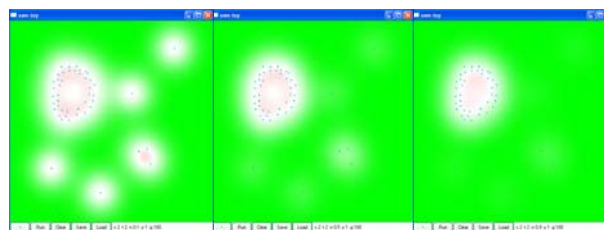
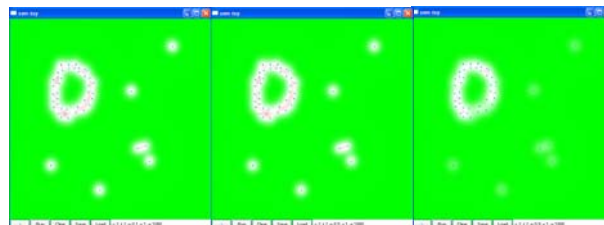
- Dans les deux formulations, si pour la solution $\rho \neq 0$, alors $v \in (0, 1]$ est
 - ◆ Une borne supérieure pour la fraction de *outliers*
 - ◆ Une borne inférieure pour la fraction de vecteurs de support
- Bornes de généralisation (formulation 2, [SPS99]) : la probabilité pour que de nouveaux exemples (tirages i.i.d. suivant la densité $p(\mathbf{x})$) soient en dehors d'une région un peu plus grande que le support déterminé ne sera pas supérieure de beaucoup à la fraction de *outliers* dans les données d'apprentissage $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$

25 septembre 2013

RCP209

29

One-class SVM : exemples

Noyau RBF avec $\gamma=100$: $v=0,1$, $v=0,5$, $v=0,9$ Noyau RBF avec $\gamma=1000$: $v=0,1$, $v=0,5$, $v=0,9$

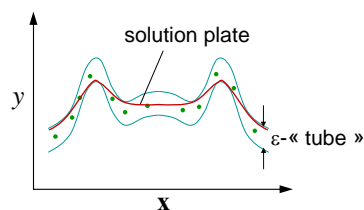
25 septembre 2013

RCP209

30

SVM pour la régression

- Données d'apprentissage $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{1 \leq i \leq n}$, $\mathbf{x} \in \mathcal{X}$, $y \in \mathbb{R}$
- En **régression ε -SV** on cherche une fonction $f : \mathcal{X} \rightarrow \mathbb{R}$, aussi « plate » que possible, et telle que $|f(\mathbf{x}_i) - y_i| \leq \varepsilon$
- On cherchera des solutions de la forme $f(\mathbf{x}) = \langle \mathbf{w}, \phi(\mathbf{x}) \rangle + b$ dans l'espace d'arrivée, la condition d'aplatissement se traduisant par la minimisation de $\|\mathbf{w}\|^2 = \langle \mathbf{w}, \mathbf{w} \rangle$



Remarque : la régression ε -SV correspond à l'utilisation de la fonction de coût **ε -insensible**

$$|\xi|_\varepsilon = \begin{cases} 0 & \text{si } |\xi| < \varepsilon \\ |\xi| - \varepsilon & \text{sinon} \end{cases}$$

25 septembre 2013

RCP209

31

Problème d'optimisation

- Comme en discrimination, on accepte quelques erreurs au-delà de ε et on introduit les « variables d'assouplissement » $\xi_i, \xi_i^* \geq 0$
- Le problème d'optimisation résultant sera

$$\begin{cases} \min_{\mathbf{w}, \xi, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \\ y_i - \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle - b \leq \varepsilon + \xi_i \\ \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0, 1 \leq i \leq n \end{cases}$$

- La constante $C > 0$ permet de choisir le point d'équilibre entre l'aplatissement de la solution et l'acceptation d'erreurs au-delà de ε

25 septembre 2013

RCP209

32

Résolution du problème d'optimisation

- Avec les multiplicateurs de Lagrange on obtient le problème dual

$$\begin{cases} \min_{\alpha} \frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) K(\mathbf{x}_i, \mathbf{x}_j) + \varepsilon \sum_{i=1}^n (\alpha_i + \alpha_i^*) - \sum_{i=1}^n y_i (\alpha_i - \alpha_i^*) \\ C \geq \alpha_i, \alpha_i^* \geq 0, \quad 1 \leq i \leq n \\ \sum_{i=1}^n (\alpha_i - \alpha_i^*) = 0 \end{cases}$$

- Comme $\mathbf{w} = \sum_{i=1}^n (\alpha_i - \alpha_i^*) \phi(\mathbf{x}_i)$, la fonction recherchée sera

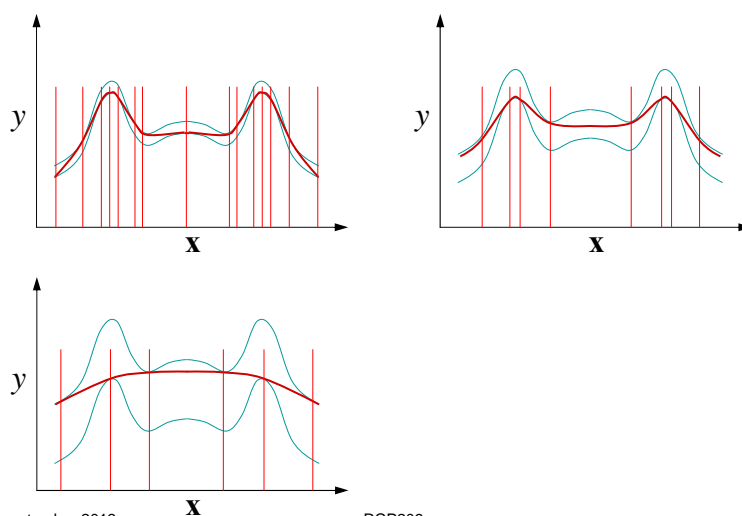
$$f(\mathbf{x}) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(\mathbf{x}, \mathbf{x}_i) + b$$

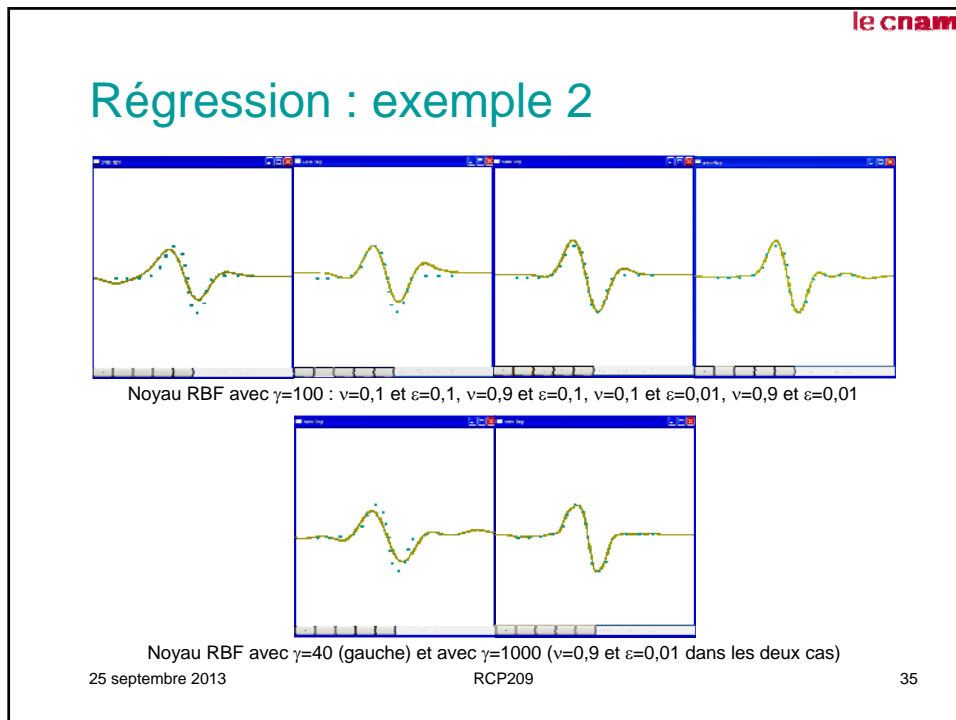
- Les conditions KKT impliquent $\alpha_i \alpha_i^* = 0$ et

$$\varepsilon - y_i + \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b \geq 0 \quad \text{et} \quad \xi_i = 0 \quad \text{si} \quad \alpha_i < C$$

$$\varepsilon - y_i + \langle \mathbf{w}, \phi(\mathbf{x}_i) \rangle + b \leq 0 \quad \text{si} \quad \alpha_i > C$$

Régression : exemple 1





le cnam

Méthodes à noyaux et analyse des données

25 septembre 2013 RCP209 36

ACP à noyaux

- L'analyse en composantes principales peut être effectuée dans l'espace d'arrivée, sur le nuage des points $\phi(\mathbf{x}_i)$

- Le problème de valeurs et vecteurs propres à résoudre est

$$\sum_{i=1}^n \langle \widehat{\phi}(\mathbf{x}_i), \mathbf{u} \rangle \widehat{\phi}(\mathbf{x}_i) = \lambda \mathbf{u}$$

- En utilisant les vecteurs centrés $\widehat{\phi}(\mathbf{x}_i) = \phi(\mathbf{x}_i) - \frac{1}{n} \sum_{j=1}^n \phi(\mathbf{x}_j)$

et en définissant la matrice $\widehat{\mathbf{K}}$ d'éléments $\langle \widehat{\phi}(\mathbf{x}_i), \widehat{\phi}(\mathbf{x}_j) \rangle$

on obtient $\widehat{\mathbf{K}} \boldsymbol{\beta} = \lambda \boldsymbol{\beta}$ où

$$\langle \widehat{\phi}(\mathbf{x}_i), \widehat{\phi}(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n} \sum_{k=1}^n K(\mathbf{x}_i, \mathbf{x}_k) - \frac{1}{n} \sum_{k=1}^n K(\mathbf{x}_j, \mathbf{x}_k) + \frac{1}{n^2} \sum_{k=1}^n \sum_{l=1}^n K(\mathbf{x}_k, \mathbf{x}_l)$$

- La projection d'un nouveau point sur un vecteur propre sera

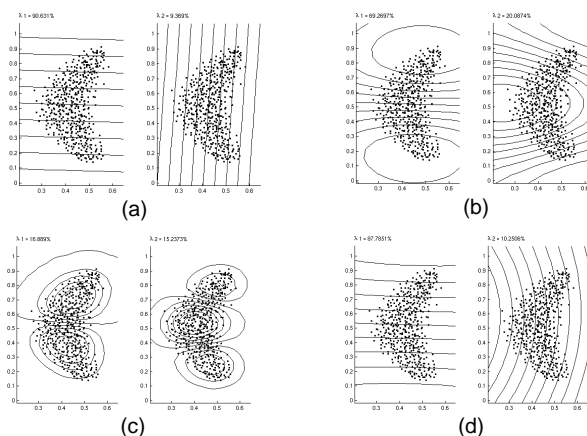
$$\langle \widehat{\phi}(\mathbf{x}), \mathbf{u}_\alpha \rangle = \sum_{j=1}^n \beta_{\alpha j} K(\mathbf{x}, \mathbf{x}_j) - \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \beta_{\alpha i} K(\mathbf{x}_i, \mathbf{x}_j)$$

25 septembre 2013

RCP209

37

ACP à noyaux : exemple



Courbes de niveau obtenues par ACP (a) linéaire, (b) à noyau RBF $\gamma=5$, (c) à noyau RBF $\gamma=100$ et (d) à noyau RBF $\gamma=0,5$ (figure de [CAB04])

25 septembre 2013

RCP209

38

AFD à noyaux

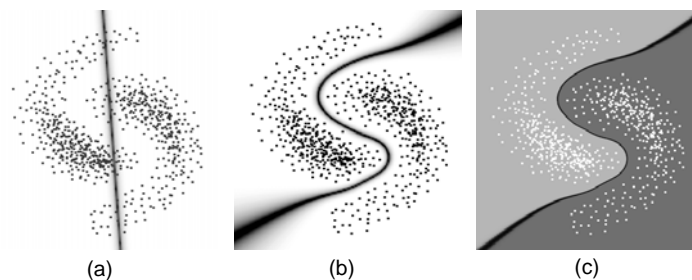
- L'analyse factorielle discriminante peut être effectuée dans l'espace d'arrivée, sur le nuage des points $\phi(x_i)$
- Une version à 2 classes a été proposée dans [MRW99] et des versions à plusieurs classes dans [BA00], [RS00]
- Si les SVM se focalisent sur les vecteurs les plus proches de la frontière (la maximisation de marge mène à une solution creuse), l'analyse discriminante à noyaux tient compte de la répartition des « masses » de points et la solution est en général moins creuse

25 septembre 2013

RCP209

39

AFD à noyaux : exemple 1



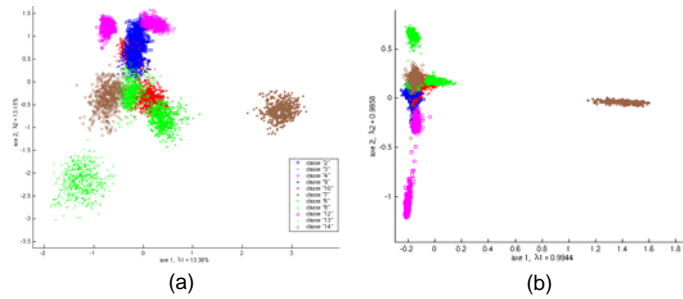
Discrimination obtenue par (a) AFD linéaire, (b) AFD à noyau angulaire et (c) machine à vecteurs support (SVM) à noyau angulaire (figure de [CAB04])

25 septembre 2013

RCP209

40

AFD à noyaux : exemple 2

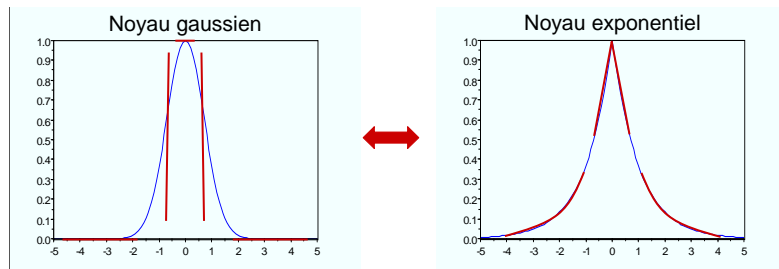


Projections des individus des données « Textures » sur les deux premiers axes factoriels discriminants pour (a) AFD linéaire et (b) AFD à noyau angulaire (figure de [CAB04])

Ingénierie des noyaux

Qu'est-ce qu'un noyau ?

- Noyau \rightsquigarrow mesure de similarité
- Définition d'une mesure de similarité $s : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$:
 - $\mathbf{x}, \mathbf{y} \in \mathcal{X} \quad s(\mathbf{x}, \mathbf{y}) \geq 0 \quad s(\mathbf{x}, \mathbf{y}) = s(\mathbf{y}, \mathbf{x})$
 - $\forall \mathbf{y}, s(\mathbf{x}, \mathbf{y}) \leq s(\mathbf{x}, \mathbf{x}) \quad s(\mathbf{x}, \mathbf{y}) = s(\mathbf{x}, \mathbf{x}) \Leftrightarrow \mathbf{x} = \mathbf{y}$



25 septembre 2013

RCP209

43

Noyaux valides

- Soit \mathcal{X} compact dans \mathbb{R}^d et $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ symétrique tels que
 - $\forall f \in L_2(\mathcal{X}), \int_{\mathcal{X}} K(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0$ (condition de Mercer)
 alors il existe un espace de Hilbert \mathcal{H} et $\phi : \mathcal{X} \rightarrow \mathcal{H}$ tels que
 - $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$
- Condition équivalente (noyau défini positif) : $\forall n \in \mathbb{N}^*$ et $\{\mathbf{x}_i\}_n \subset \mathcal{X}$, la matrice de Gram $\mathbf{K}_n : \{\mathbf{K}_n\}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j), 1 \leq i, j \leq n$, est définie positive, c'est à dire $\forall \mathbf{c}_n \in \mathbb{R}^n$ différent du vecteur nul de \mathbb{R}^n on a $\mathbf{c}_n^T \mathbf{K}_n \mathbf{c}_n > 0$
- Un noyau valide garantit donc l'existence de \mathcal{H} et peut s'exprimer comme un produit scalaire dans \mathcal{H} ; aussi, il garantit la convexité du problème d'optimisation quadratique sous contraintes

25 septembre 2013

RCP209

44

Noyaux valides (2)

- Noyau **conditionnellement défini positif (cdp)** : $\forall n \in \mathbb{N}^*, \{ \mathbf{x}_i \}_n \subset \mathcal{X}$, la matrice de Gram $\mathbf{K}_n : \{ \mathbf{K}_n \}_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$, $1 \leq i, j \leq n$, est conditionnellement définie positive, c'est à dire $\forall \mathbf{c}_n \in \mathbb{R}^n$ différent du vecteur nul de \mathbb{R}^n et tel que $\sum_{i=1}^n c_i = 0$, on a $\mathbf{c}_n^T \mathbf{K}_n \mathbf{c}_n > 0$
- Étant donné un noyau symétrique K (un noyau **cdp** est symétrique), il existe un espace vectoriel \mathcal{V} , une transformation $\phi : \mathcal{X} \rightarrow \mathcal{V}$ et une forme bilinéaire $Q : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$ tels que $K(\mathbf{x}, \mathbf{y}) = Q[\phi(\mathbf{x}), \phi(\mathbf{y})]$ (si K n'est pas défini positif, Q ne sera pas un produit scalaire) [Sch01]
- Un noyau **cdp** peut être utilisé pour les SVM en discrimination car les contraintes du problème d'optimisation quadratique incluent la condition $\sum_{i=1}^n \alpha_i y_i = 0$ ($c_i = \alpha_i y_i$)

25 septembre 2013

RCP209

45

Construire des noyaux définis positifs

- **Construction directe** : définition de \mathcal{H} , $\phi : \mathcal{X} \rightarrow \mathcal{H}$ et ensuite du noyau $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ par $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$
- Soit $f : \mathcal{X} \rightarrow \mathbb{R}$, alors $K(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) \cdot f(\mathbf{y})$ ($K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$, \mathcal{X} compact dans \mathbb{R}) est défini positif (*conformal kernel*)
 - ◆ Exemples (ces noyaux ne peuvent pas être interprétés comme des mesures de similarité) :
 1. $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = x$, $K(x, y) = x \cdot y$ (le noyau linéaire)
 2. $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = e^x$, $K(x, y) = e^{x+y}$

25 septembre 2013

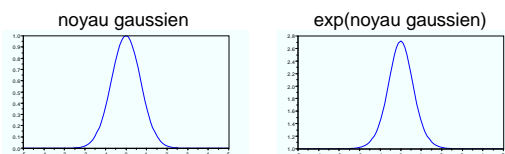
RCP209

46

Transformer des noyaux

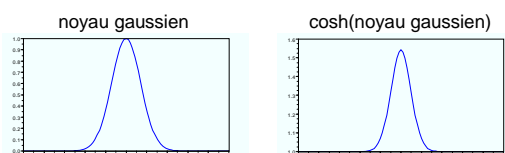
- Si $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ est défini positif alors $\exp(K)$ est défini positif

◆ Exemple :



- Si $K : \mathcal{X} \times \mathcal{X} \rightarrow [-1, 1]$ est défini positif alors $\cosh(K)$ est défini positif

◆ Exemple :



25 septembre 2013

RCP209

47

Combiner des noyaux définis positifs

- Si $K_1, K_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ sont définis positifs et $\alpha_1, \alpha_2 \in \mathbb{R}^+$, alors sont également définis positifs les noyaux suivants

◆ Combinaison linéaire : $K = \alpha_1 K_1 + \alpha_2 K_2$

◆ Produit simple : $K = K_1 \cdot K_2$

avec $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

- Si $K_1 : \mathcal{X}_1 \times \mathcal{X}_1 \rightarrow \mathbb{R}$ et $K_2 : \mathcal{X}_2 \times \mathcal{X}_2 \rightarrow \mathbb{R}$ sont définis positifs, alors sont également définis positifs les noyaux

◆ Somme directe : $K_1 \oplus K_2 = K_1 + K_2$

◆ Produit tensoriel : $K_1 \otimes K_2 = K_1 \cdot K_2$

avec $K_1 \oplus K_2, K_1 \otimes K_2 : (\mathcal{X}_1 \times \mathcal{X}_2) \times (\mathcal{X}_1 \times \mathcal{X}_2) \rightarrow \mathbb{R}$

→ construction de **noyaux hybrides**

25 septembre 2013

RCP209

48

Qu'est-ce qu'un « bon » noyau ?

- Propriétés d'un « bon » noyau (suivant [Gär03]) :
 - ◆ Préalable : définition de « concepts » $c : \mathcal{X} \rightarrow \{0, 1\}$
 - 1. « Complétude » : dans quelle mesure la connaissance incorporée dans le noyau est suffisante pour résoudre le problème ($K(\mathbf{x}, \cdot) = K(\mathbf{y}, \cdot) \Rightarrow c(\mathbf{x}) = c(\mathbf{y})$)
 - 2. « Correctitude » : dans quelle mesure la sémantique du problème est respectée par le noyau (l'outil qui exploite le noyau, SVM dans notre cas, permet de trouver une solution pour chaque concept du problème)
 - 3. « Justesse » : dans quelle mesure le noyau reflète bien la mesure de similarité qui correspond aux « concepts » du problème (des bornes polynomiales de généralisation peuvent être obtenues pour l'outil qui exploite le noyau)

25 septembre 2013

RCP209

49

Exemple : noyau pour ensembles

- Nature des données (exemple issu de [Bou05]) : ensembles de descripteurs locaux d'images
- Définition de l'espace \mathcal{X} : ensemble des parties finies mais de cardinalité variable de \mathbb{R}^d : $\mathcal{X} = \mathcal{P}_f(\mathbb{R}^d)$
- Objectif : évaluer la similarité entre ensembles de vecteurs $\mathcal{E}, \mathcal{E}' \in \mathcal{P}_f(\mathbb{R}^d)$ à travers les proximités entre vecteurs similaires de $\mathcal{E}, \mathcal{E}' \rightarrow$ on décrit ici le **noyau d'appariement intermédiaire** (*intermediate matching kernel*, [Bou05])
- Problème : le noyau d'appariement direct,

$$K_{\mathcal{A}}(\mathcal{E}, \mathcal{E}') = (1/2) \left[\sum_{\mathbf{x}_i \in \mathcal{E}} \max_{\mathbf{x}'_j \in \mathcal{E}'} K(\mathbf{x}_i, \mathbf{x}'_j) + \sum_{\mathbf{x}'_j \in \mathcal{E}'} \max_{\mathbf{x}_i \in \mathcal{E}} K(\mathbf{x}_i, \mathbf{x}'_j) \right]$$

où $K(\mathbf{x}_i, \mathbf{x}'_j)$ est un noyau classique entre les vecteurs $\mathbf{x}_i, \mathbf{x}'_j$, n'est pas défini positif !

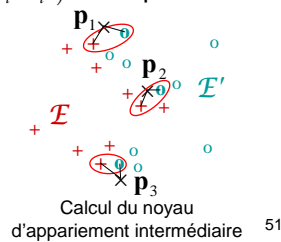
25 septembre 2013

RCP209

50

Exemple : noyau pour ensembles (2)

- Principe du noyau d'appariement intermédiaire : faire l'appariement par rapport à des vecteurs pivots, fixés pour un ensemble d'apprentissage donné (ne dépendant donc pas de $\mathcal{E}, \mathcal{E}'$)
- Soient m vecteurs pivots $\mathbf{p}_1, \dots, \mathbf{p}_m \in \mathbb{R}^d$, pour chaque vecteur \mathbf{p}_i on définit une fonction $\psi_i : \mathcal{X} \rightarrow \mathbb{R}^d, \psi_i(\mathcal{E}) = \mathbf{x}_i^* = \arg \min_{\mathbf{x} \in \mathcal{E}} \|\mathbf{x} - \mathbf{p}_i\|$
- Le noyau d'appariement intermédiaire construit à partir des pivots $\mathbf{p}_1, \dots, \mathbf{p}_m$ est défini par $K_{\mathcal{M}}(\mathcal{E}, \mathcal{E}') = \sum_{l=1}^m K(\mathbf{x}_l^*, \mathbf{x}_l'^*)$ et on peut facilement montrer qu'il est positif défini
- Choix possible des vecteurs pivots : prototypes des groupes de vecteurs obtenus par classification automatique des données d'apprentissage



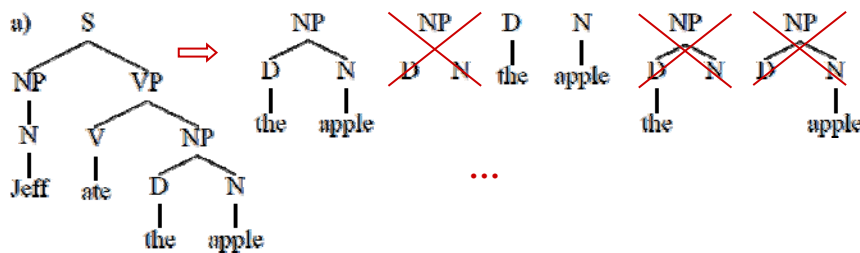
25 septembre 2013

RCP209

Calcul du noyau d'appariement intermédiaire 51

Exemple : noyau pour arbres

- Nature des données (exemple issu de [CD02]) : arbres de *parsing* employés dans le traitement automatique des langues
- Définition de l'espace \mathcal{X} : ensemble des arbres étiquetés et ordonnés, avec des étiquettes dans un ensemble E
- Sous-arbre : sous-graphe connexe d'un arbre qui, pour un nœud de l'arbre, contient soit tous ses fils, soit aucun



25 septembre 2013

RCP209

52

Exemple : noyau pour arbres (2)

- On considère une énumération de tous les sous-arbres possibles (nombre fini)
- On note par $h_l(T)$ le nombre d'occurrences du sous-arbre l dans l'arbre T
- L'arbre T sera représenté par $\phi(T) = [h_1(T) h_2(T) \dots h_l(T) \dots]$
- Pour deux arbres $T_1, T_2 \in \mathcal{X}$ on définit la valeur du noyau comme étant le produit scalaire entre ces deux représentations, $K(\mathbf{x}, \mathbf{y}) = \langle \phi(\mathbf{x}), \phi(\mathbf{y}) \rangle$ ou, plus précisément, $K(T_1, T_2) = \sum_l h_l(T_1) h_l(T_2)$
- Soit $\mathcal{N}_1, \mathcal{N}_2$ les ensembles des nœuds des arbres T_1, T_2 et $S(v_1, v_2)$, $v_1 \in \mathcal{N}_1, v_2 \in \mathcal{N}_2$ le nombre de sous-arbres isomorphes ayant pour racine v_1 et respectivement v_2
- Alors $K(T_1, T_2) = \sum_{v_1 \in \mathcal{N}_1, v_2 \in \mathcal{N}_2} S(v_1, v_2)$, calcul récursif $O(|\mathcal{N}_1| \cdot |\mathcal{N}_2|)$

25 septembre 2013

RCP209

53

Exemples d'application

25 septembre 2013

RCP209

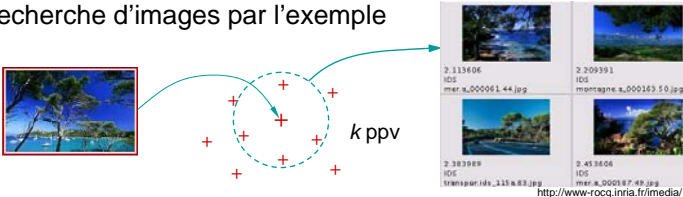
54

le cnam

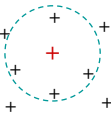
SVM pour le contrôle de pertinence

- Contrôle de pertinence (*relevance feedback*) : tenir compte du feedback de l'utilisateur dans la recherche itérative par le contenu

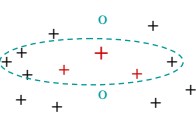
1. Recherche d'images par l'exemple


2. Recherche itérative avec contrôle de pertinence

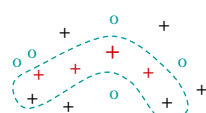
1



2



3



25 septembre 2013
RCP209
55

le cnam

Composantes du mécanisme

1. **Learner** : à partir de l'information disponible (notamment des exemples positifs et/ou négatifs), estimer l'ensemble d'images visé
2. **Sélecteur** : à partir de l'estimation produite par le learner, choisir les images que l'utilisateur doit marquer lors de l'itération suivante
3. Utilisateur : fournir à chaque itération le retour pour les images choisies par le sélecteur
 - Les évaluations sont souvent faites à l'aide d'une vérité terrain, en émulant l'utilisateur

25 septembre 2013
RCP209
56

Difficultés pour l'apprentissage

- **Très peu d'exemples** étiquetés : leur nombre est souvent inférieur au nombre de dimensions de l'espace de description !
- **Déséquilibre** important entre le nombre d'exemples positifs et le nombre d'exemples négatifs
- **Forme** potentiellement **complexe** de l'ensemble d'images visé, qui peut même présenter **plusieurs modes distants** dans l'espace de description
- L'interactivité exige un **temps de réponse très court**, à la fois pour le *learner* et pour le sélecteur

25 septembre 2013

RCP209

57

Sélecteur : objectifs et critères

- Objectifs
 1. Retourner un maximum d'images pertinentes à l'utilisateur
 2. Maximiser le transfert d'information **utilisateur → système**
- Critères de sélection
 - ◆ « **Les plus positives** » (MP) : retourner les images les plus pertinentes suivant l'estimation actuelle faite par le *learner* – critère classique le plus utilisé
 - ◆ « **Les plus informatives** » (MI) : retourner les images qui permettent à l'utilisateur de fournir un maximum d'information sur sa cible → **minimiser le nombre d'exemples**

25 septembre 2013

RCP209

58

Critère « les plus informatives »

- Composantes **complémentaires** du critère MI
 1. Ambiguïté élevée (de chaque image sélectionnée) par rapport à l'estimation courante faite par le *learner*
 - ◆ Comme critère individuel : « les plus ambiguës »
 2. Faible redondance de l'ensemble des s images retournées
- Un critère « les plus informatives » pour SVM [FCB04]
 1. Présélectionner les $t > s$ images pour lesquelles les valeurs de la fonction de décision SVM sont les plus proches de 0 (images les plus ambiguës)
 2. Choisir itérativement les s images pour lesquelles

$$\mathbf{x}_j = \arg \min_{\mathbf{x} \in \text{présélection}} \max_i K(\mathbf{x}, \mathbf{x}_i)$$

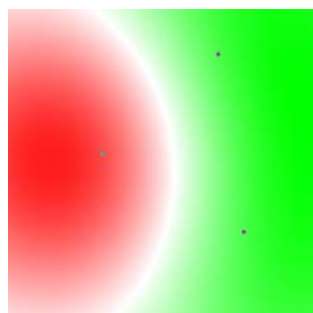
25 septembre 2013

RCP209

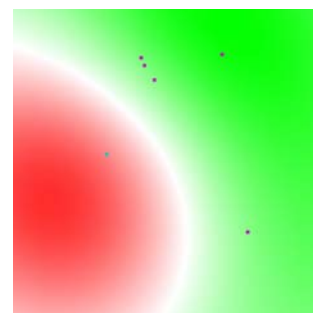
59

Plus ambiguës : illustration

- Les s images les plus proches de la frontière peuvent être redondantes



Avant sélection

Après sélection, *feedback*, estimation

25 septembre 2013

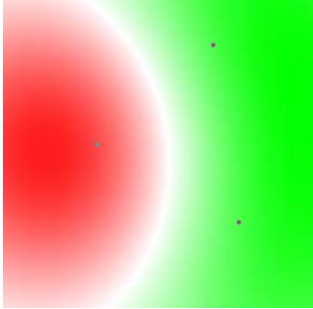
RCP209

60

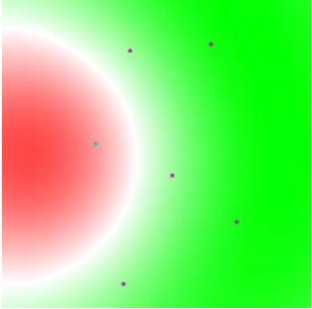
le cnam

Plus informatives : illustration

- Le critère conjoint minimise également la redondance
 ⇒ focalisation plus rapide sur les images recherchées



Avant sélection



Après sélection, *feedback*, estimation

25 septembre 2013
RCP209
61

le cnam


Contrôle de pertinence : exemple (3)

Objectif : retrouver des portraits

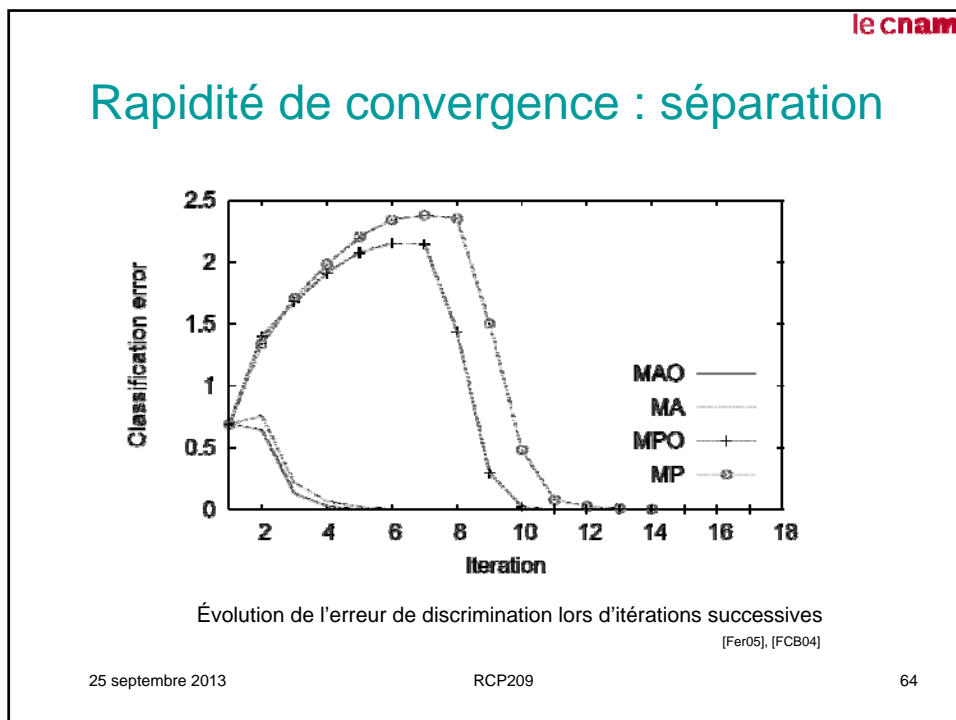
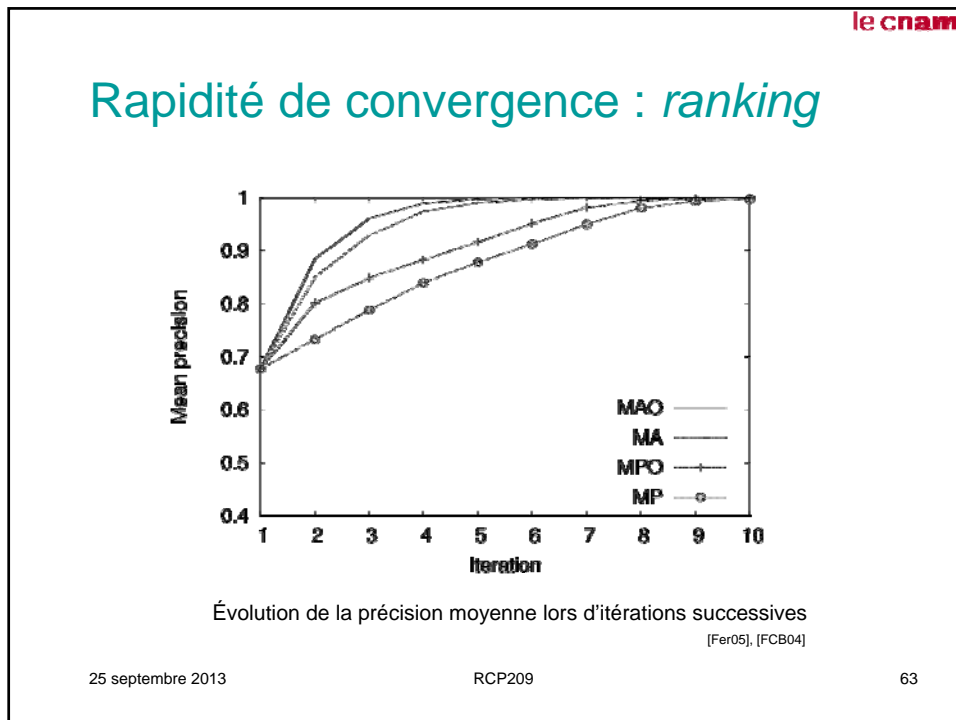
Base de 7500 images, dont 110 portraits

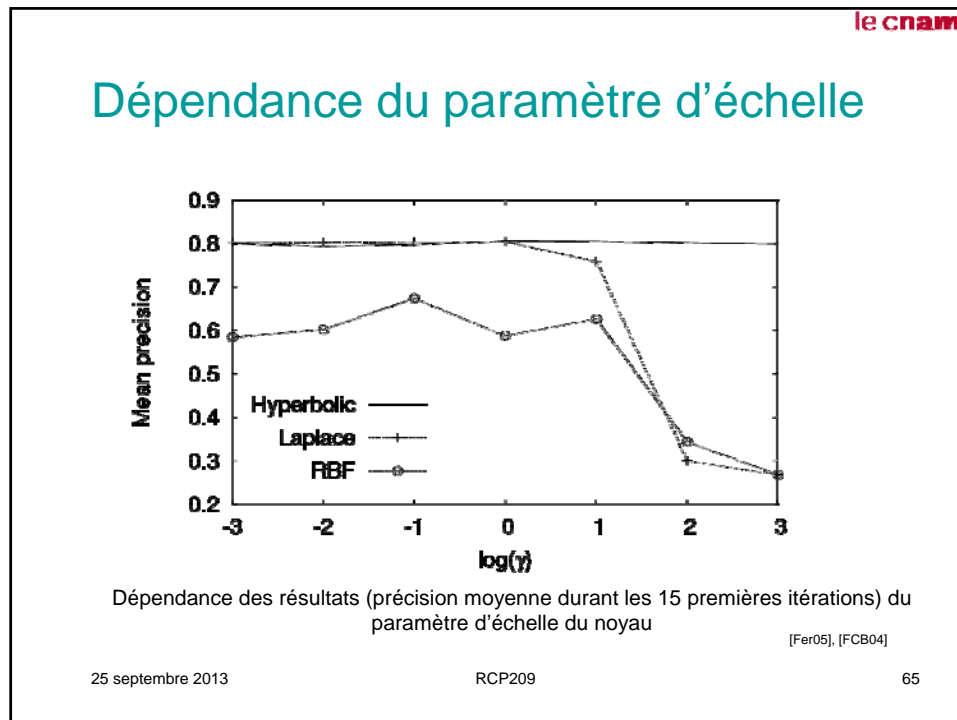
Disponible : description globale (couleur, texture, forme)

Première page de résultats après 4 itérations →



25 septembre 2013
RCP209
Voir [FCB04] et <http://www-rocq.inria.fr/imedia62>





le cnam

Application : conclusion

- Avantages des SVM pour le contrôle de pertinence
 - ◆ La fonction de décision associée permet à la fois la définition d'une frontière et le classement des images
 - ◆ Avec un large choix des noyaux, les SVM permettent une grande liberté dans la forme des classes (avec un contrôle par la régularisation)
 - ◆ D'autres sources d'information (en dehors des exemples) permettent de définir des noyaux appropriés
 - ◆ Apprentissage très rapide avec le nombre relativement limité d'exemples fournis par le contrôle de pertinence
 - ◆ Moindre sensibilité au déséquilibre entre exemples positifs et négatifs
- Inconvénients
 - ◆ Par rapport aux noyaux de Parzen, absence de caractère incrémental (dans la formulation de base) et donc étape de sélection plus coûteuse

25 septembre 2013 RCP209 66

Références

- [BA00] G. Baudat, Anouar F., Generalized discriminant analysis using a kernel approach. *Neural Computation*, 12: 2385-2404, 2000.
- [Bou05] S. Boughorbel, Noyaux pour la classification d'images par les Machines à Vecteurs de Support. *Thèse de doctorat*, Université d'Orsay, juillet 2005.
- [Bur98] C.J.C. Burges. A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery 2*: 121-167, 1998.
- [CAB04] M. Crucianu, J.-P. Asselin de Beauville, R. Boné, *Méthodes factorielles pour l'analyse des données : méthodes linéaires et extensions non-linéaires*. Hermès, 2004, 288 p.
- [CD02] M. Collins, N. Duffy. Convolution kernels for natural language. In T.G. Dietterich, S. Becker and Z. Ghahramani (Eds.) *Advances in Neural Information Processing Systems* Vol. 15. MIT Press, 2002.
- [Fer05] M. Ferecatu. Image retrieval with active relevance feedback using both visual and keyword-based descriptors, *Thèse de doctorat*, Université de Versailles Saint-Quentin-en-Yvelines, juillet 2005.

25 septembre 2013

RCP209

67

Références

- [FCB04] Ferecatu, M., Crucianu, M., Boujemaa, N. Retrieval of Difficult Image Classes Using SVM-Based Relevance Feedback, *Proceedings of the 6th ACM SIGMM International Workshop on Multimedia Information Retrieval*, 2004, New York, pp. 23-30.
- [Gär03] Gärtner, T. A survey of kernels for structured data. *SIGKDD Explorer Newsletter* 5: 49-58, 2003.
- [MRW99] S. Mika, Rätsch G., Weston J., Schölkopf B., Muller K.-R., Fischer discriminant analysis with kernels, *Proceedings of the IXth international workshop on Neural Networks for Signal Processing*, pp. 41-48, IEEE Press, 1999.
- [MMR01] K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, B. Schölkopf. An introduction to kernel-based learning methods. *IEEE Transactions on Neural Networks* 12: 181-202, 2001.
- [RS00] V. Roth, Steinhage V., Nonlinear discriminant analysis using kernel functions, *Advances in Neural Information Processing Systems*, vol. 12, pp. 568-574, MIT Press, 2000.

25 septembre 2013

RCP209

68

Références

- [Sch01] Schölkopf, B. The kernel trick for distances. *Advances in Neural Information Processing Systems* 13, 301-307. (Eds.) Leen, T. K., T. G. Dietterich, V. Tresp, MIT Press, Cambridge, MA, USA, 2001.
- [SPS99] B. Schölkopf, J. Platt, A. Smola, J. Shawe-Taylor, R. J. Williamson. Estimating the support of a high-dimensional distribution. *Technical Report* 87, Microsoft Research, Redmond, WA, 1999.
- [SSM96] B. Schölkopf, Smola A., Müller K.-R., Nonlinear component analysis as a kernel eigenvalue problem, *Technical Report* 44, Max-Planck-Institut für biologische Kybernetik, Tübingen, Germany, 1996.
- [SS02] B. Schölkopf, A. Smola. *Learning with Kernels*. MIT Press, 2002.
- [SS04] A. Smola, B. Schölkopf. A tutorial on support vector regression. *Statistics and Computing* 14: 199-222, 2004.
- [TD99] D. Tax, R. Duin. Data domain description by support vectors. *Proceedings ESANN'99*, pp. 251-256. D-Facto Press, Brussels. 1999.
- [Vap98] V. Vapnik. *Statistical Learning Theory*. John Wiley, New York. 1998.