



Apprentissage et noyaux séparateur à vaste marge (SVM)

Pour quoi faire ?

9 Novembre 2006

Stéphane Canu

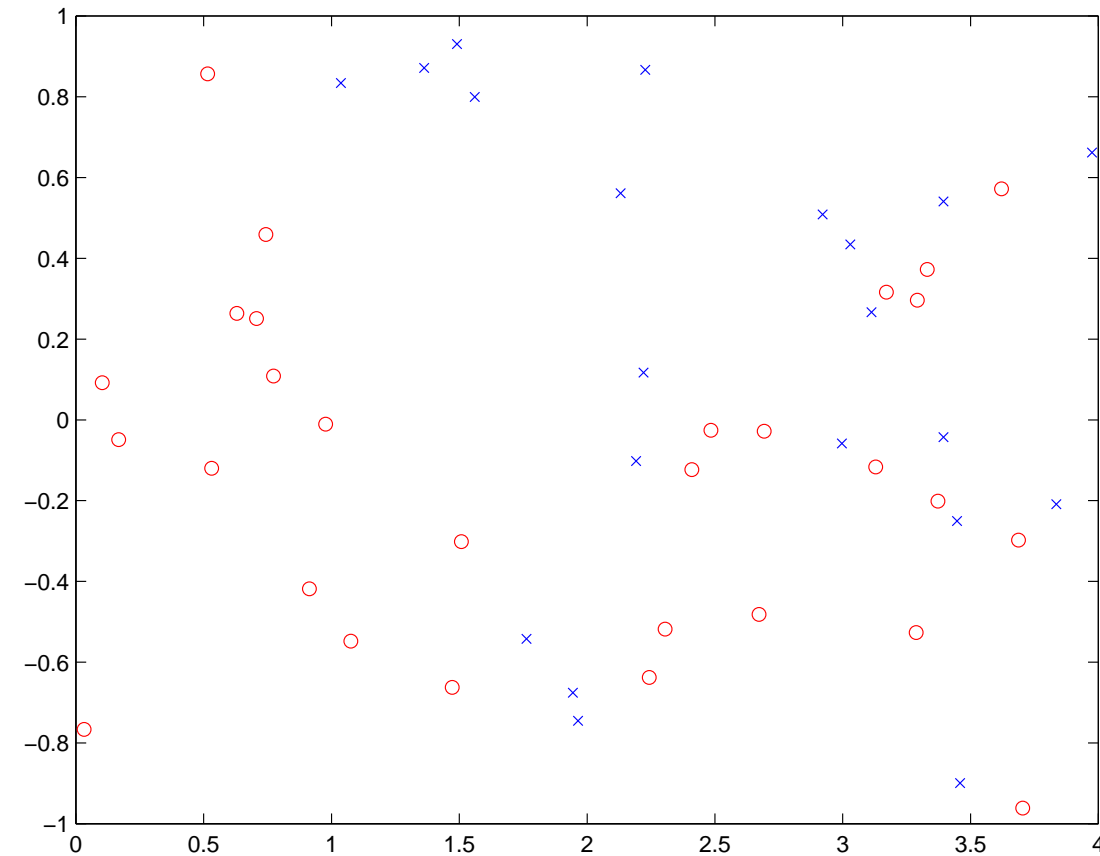
`stephane.canu@insa-rouen.fr`

`asi.insa-rouen.fr/~scanu`

INSA Rouen - Département ASI

Laboratoire LITIS

A la recherche d'une règle de décision universelle



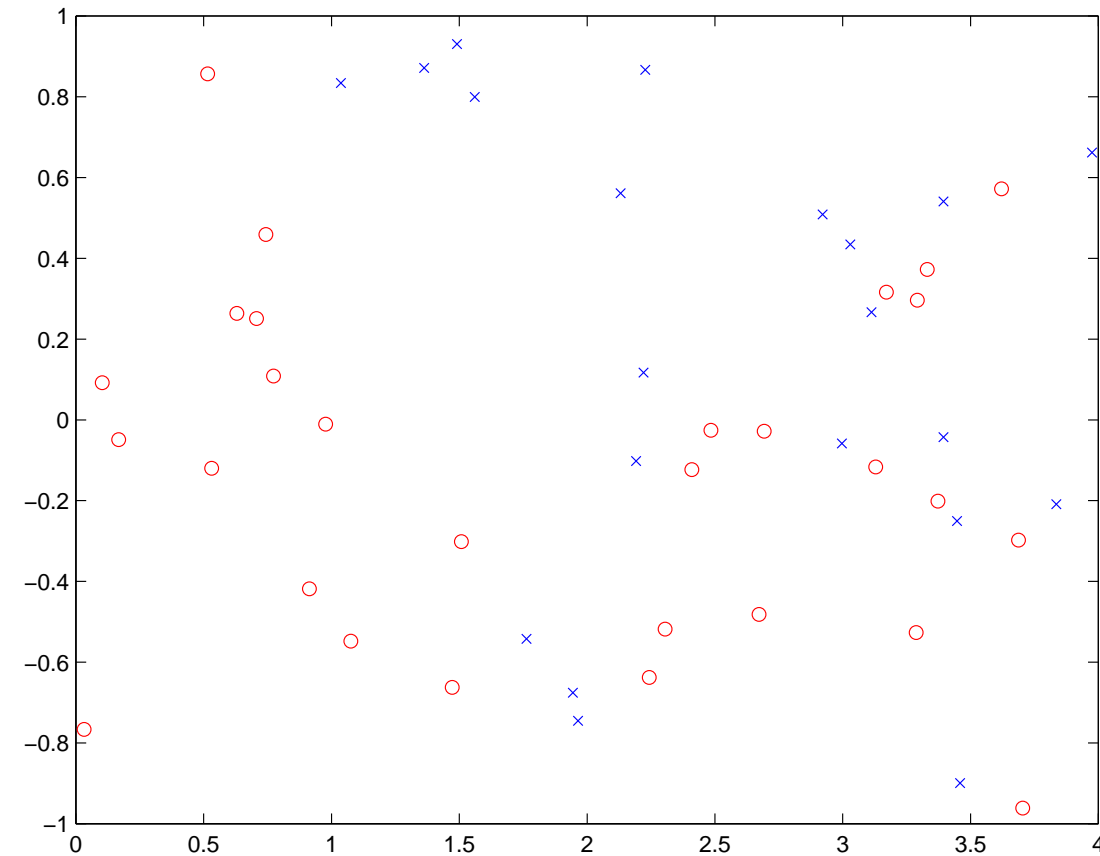
on cherche un algorithme \mathcal{A}
capable de résoudre
tous les problèmes

l'échantillon $(x_i, y_i)_{i=1,n}$

$$\underbrace{\mathbb{P}(\text{err}(f, x_i, y_i))}_{\text{erreur de } f} \xrightarrow{n \rightarrow \infty} \underbrace{\mathbb{P}_{\mathbf{b}}(\text{err})}_{\text{erreur de bayes}}$$

Tracez la frontière de décision entre ces deux classes ?

A la recherche d'une règle de décision universelle



on cherche un algorithme \mathcal{A}
capable de résoudre
tous les problèmes

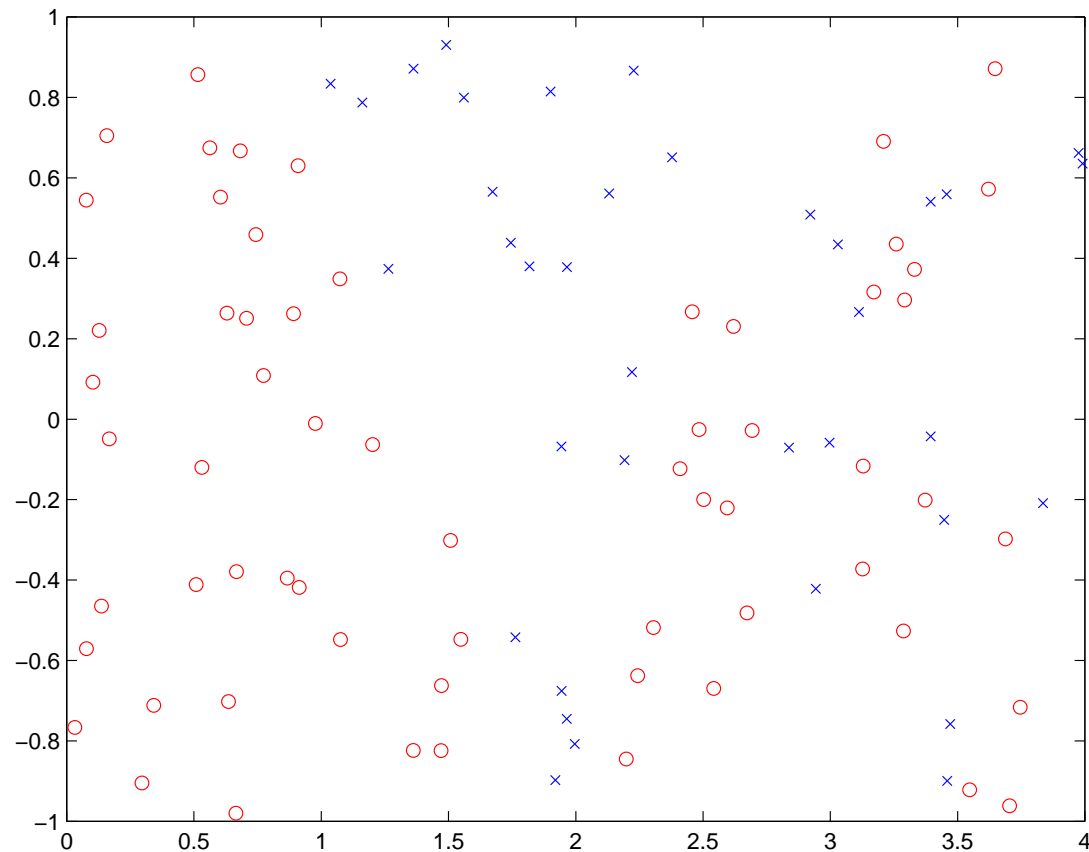
l'échantillon $(x_i, y_i)_{i=1,n}$

$$\underbrace{\mathbb{P}(\text{err}(f, x_i, y_i))}_{\text{erreur de } f} \xrightarrow{n \rightarrow \infty} \underbrace{\mathbb{P}_{\mathbf{b}}(\text{err})}_{\text{erreur de bayes}}$$

Tracez la frontière de décision entre ces deux classes ?

Universelle : pour tous les problèmes

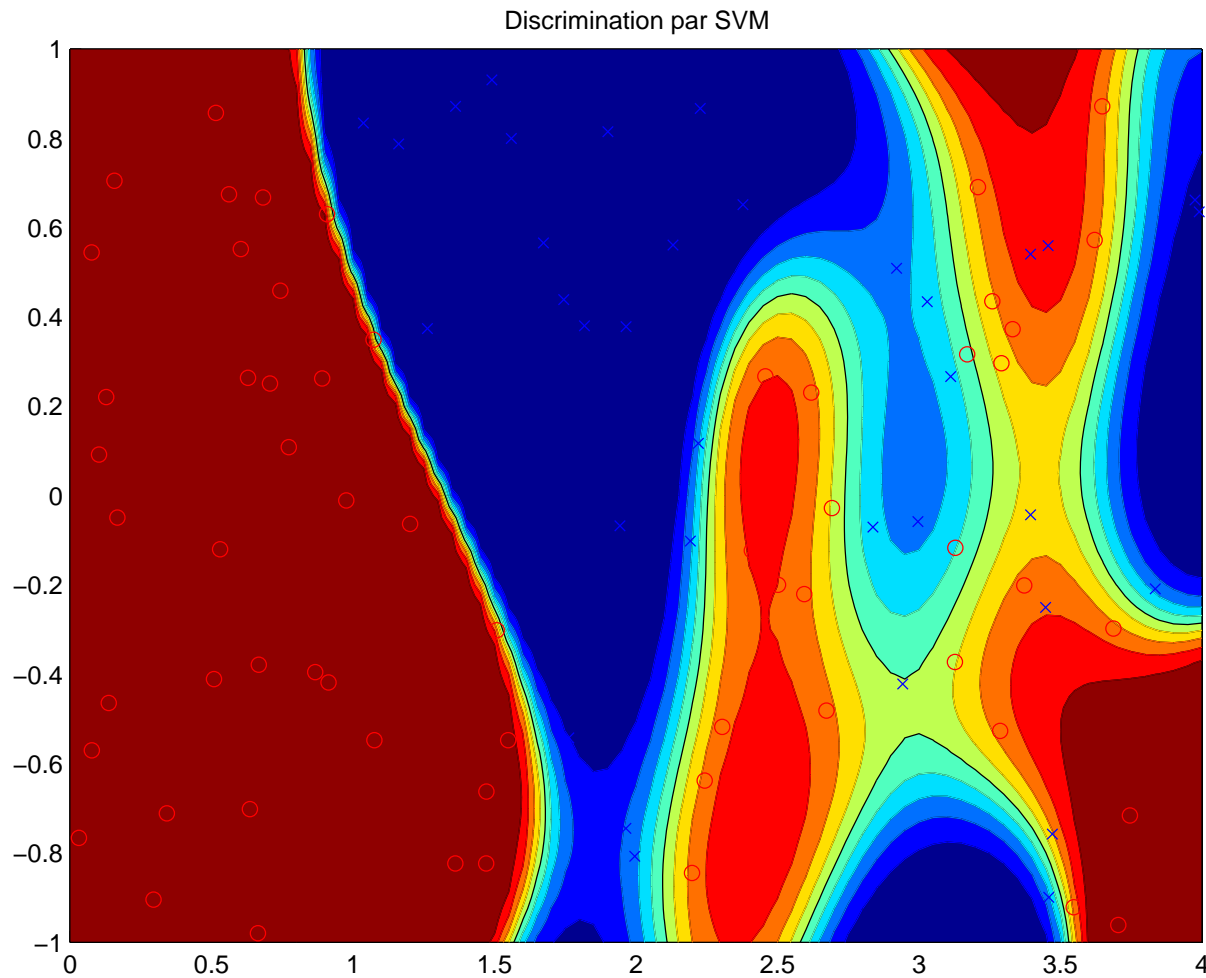
Introduction



c'est plus facile...
...avec un peu plus de points

Tracez la frontière de décision entre ces deux classes ?

Introduction : le plan

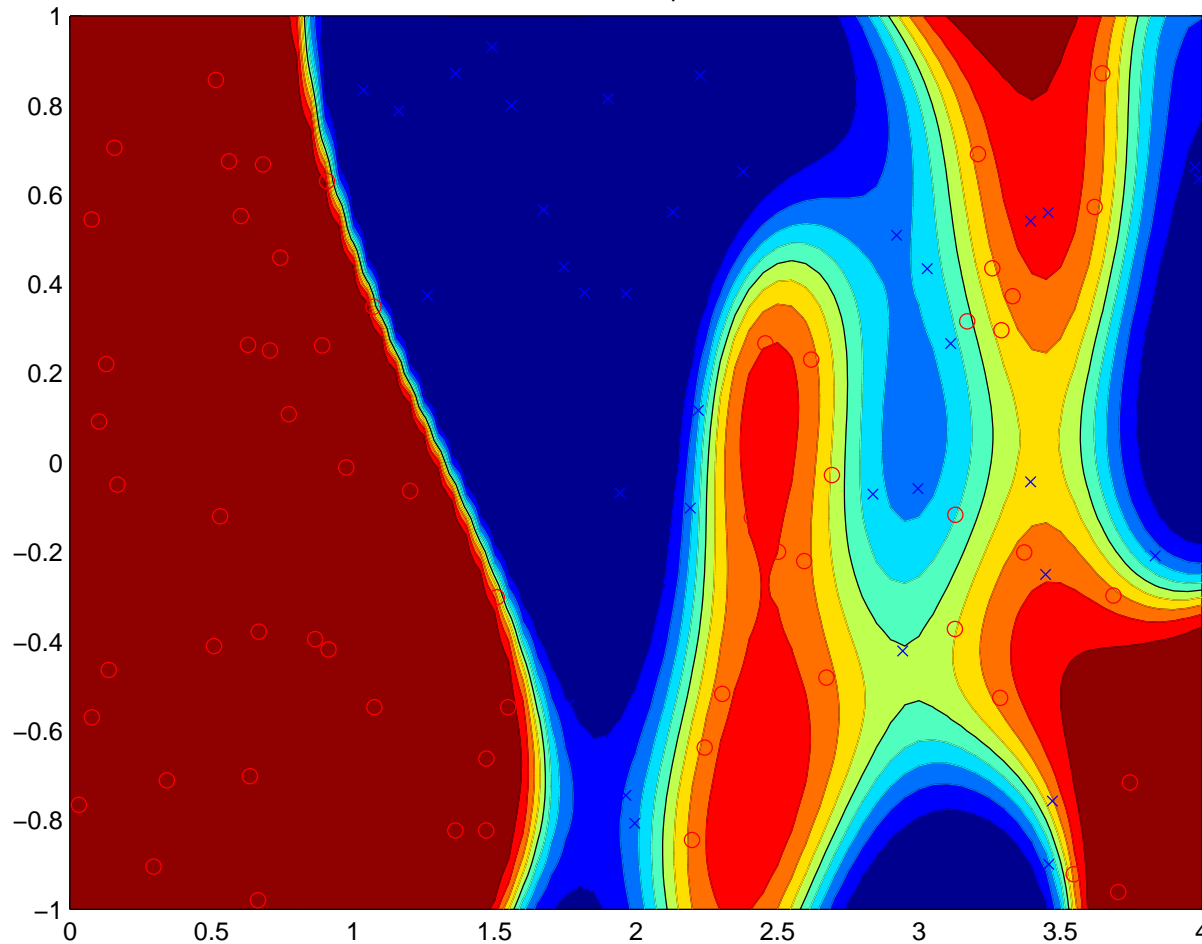


Une solution
⇒ Quels critères ?

■ (1) Fidélité

Introduction : le plan

Discrimination par SVM



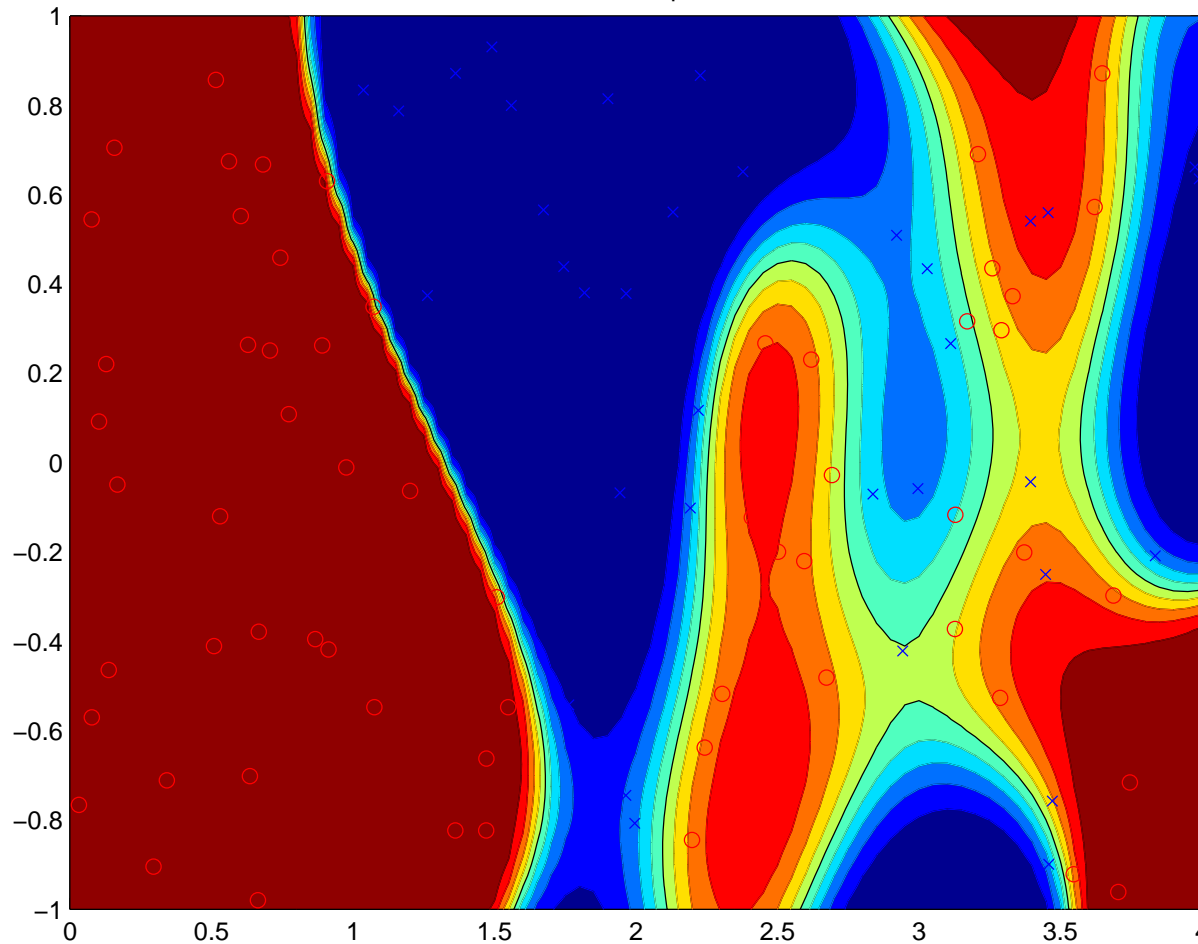
Une solution
⇒ Quels critères ?

- (1) Fidélité
- (2) Régularité

Théorème [Vapnik, 1979] : $\underbrace{\mathbb{P}(\text{err}[f])}_{\text{erreur}} \leq \underbrace{R_{\text{emp}}[f]}_{\text{fidélité}} + \underbrace{\sqrt{\frac{h(\log(2n/h) + 1) - \log(\eta/4)}{n}}}_{\text{régularité}}$

Introduction : le plan

Discrimination par SVM



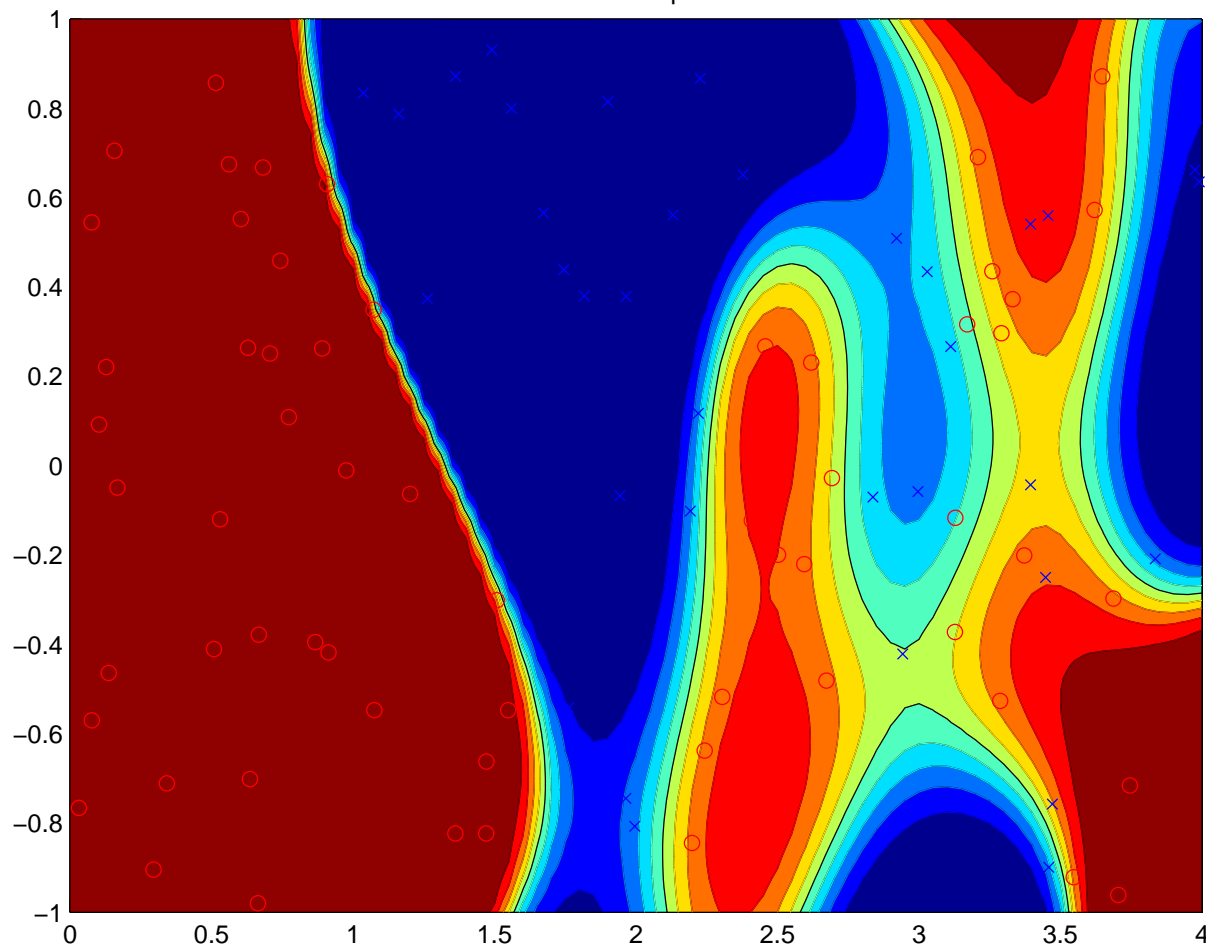
Une solution
⇒ Quels critères ?

- (1) Fidélité
- (2) Régularité
- (3) Décision locale

Théorème [Vapnik, 1979] : $\underbrace{\mathbb{P}(\text{err}[f])}_{\text{erreur}} \leq \underbrace{R_{\text{emp}}[f]}_{\text{fidélité}} + \underbrace{\sqrt{\frac{h(\log(2n/h) + 1) - \log(\eta/4)}{n}}}_{\text{régularité}}$

Introduction : le plan

Discrimination par SVM

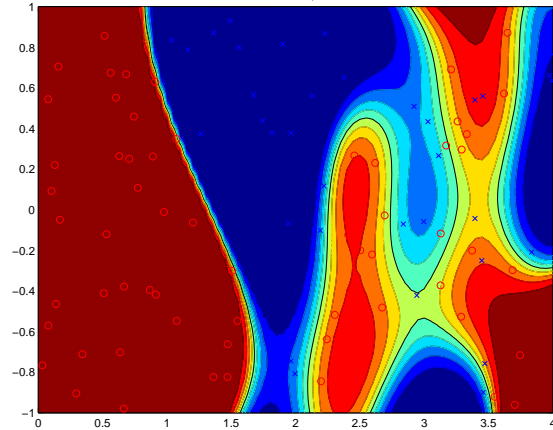


Une solution
⇒ Quels critères ?

- (1) Fidélité
- (2) Régularité
- (3) Décision locale
- (4) Points frontière

Théorème [Vapnik, 1979] : $\underbrace{\mathbb{P}(\text{err}[f])}_{\text{erreur}} \leq \underbrace{R_{\text{emp}}[f]}_{\text{fidélité}} + \underbrace{\sqrt{\frac{h(\log(2n/h) + 1) - \log(\eta/4)}{n}}}_{\text{régularité}}$

Discrimination par SVM



l'échantillon $(x_i, y_i)_{i=1,n}$

$y_i \in \{-1, 1\}$ (codage -1/1)

la fonction de décision : $\text{signe}(f(x_i))$

(f fonction de discrimination)

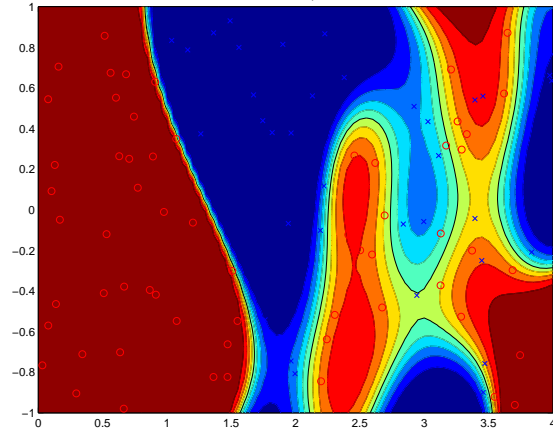
$\mathcal{H} = \{x \mid f(x) = 0\}$: frontière de décision.

Bien classer tout le monde :

$$\text{signe}(f(x_i)) = y_i \quad i = 1, n$$

(1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Discrimination par SVM



l'échantillon $(x_i, y_i)_{i=1,n}$

$y_i \in \{-1, 1\}$ (codage -1/1)

la fonction de décision : $\text{signe}(f(x_i))$

(f fonction de discrimination)

$\mathcal{H} = \{x \mid f(x) = 0\}$: frontière de décision.

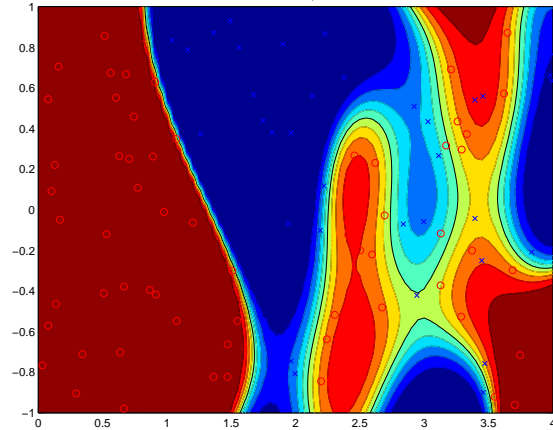
Bien classer tout le monde :

$\text{signe}(f(x_i)) = y_i \quad i = 1, n \quad \text{critère non dérivable}$

$f(x_i)y_i \geq 0 \quad i = 1, n$

(1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Discrimination par SVM



l'échantillon $(x_i, y_i)_{i=1, n}$

$y_i \in \{-1, 1\}$ (codage -1/1)

la fonction de décision : $\text{signe}(f(x_i))$

(f fonction de discrimination)

$\mathcal{H} = \{x \mid f(x) = 0\}$: frontière de décision.

Bien classer tout le monde :

$\text{signe}(f(x_i)) = y_i \quad i = 1, n \quad \text{critère non dérivable}$

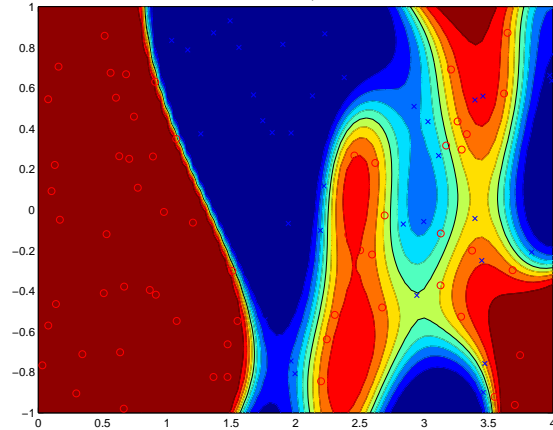
$f(x_i)y_i \geq 0 \quad i = 1, n \quad \text{solution triviale } f = 0$

$$f(x_i)y_i \geq k \quad k > 0, i = 1, n$$

(1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Fidélité

Discrimination par SVM



l'échantillon $(x_i, y_i)_{i=1, n}$

$y_i \in \{-1, 1\}$ (codage -1/1)

la fonction de décision : $\text{signe}(f(x_i))$

(f fonction de discrimination)

$\mathcal{H} = \{x \mid f(x) = 0\}$: frontière de décision.

Bien classer tout le monde :

$$\text{signe}(f(x_i)) = y_i \quad i = 1, n \quad \text{critère non dérivable}$$

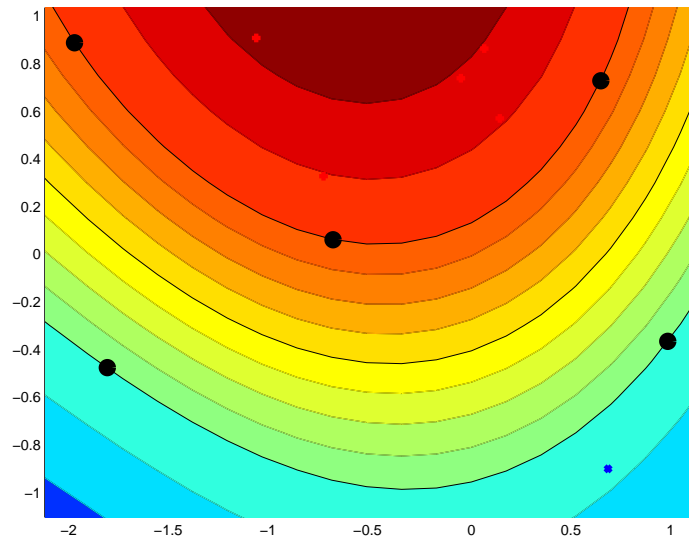
$$f(x_i)y_i \geq 0 \quad i = 1, n \quad \text{solution triviale } f = 0$$

$$f(x_i)y_i \geq k \quad k > 0, i = 1, n$$

Marge

- (1) **Fidélité** - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Fidélité et marge



$$f(x_i)y_i > k \quad k > 0, \quad i = 1, n$$

$\mathcal{H} = \{x \mid f(x) = 0\}$: frontière

Marge :

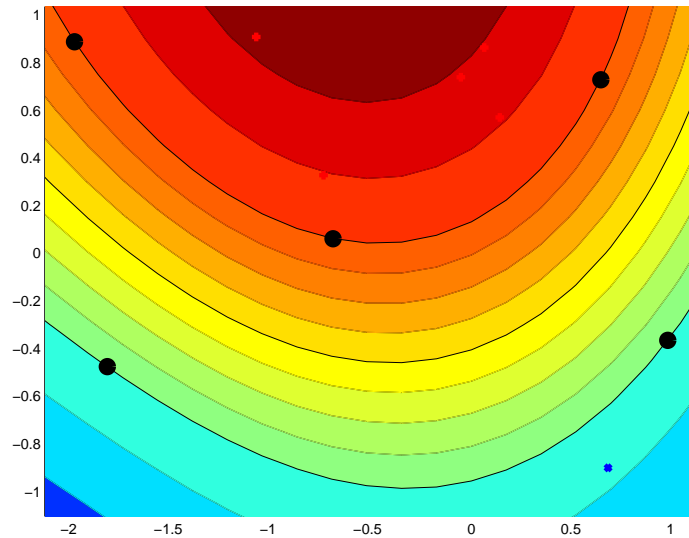
$$\min_{i=1,n} d(\mathcal{H}, x_i) = \min_{i=1,n} \max(1 - f(x_i)y_i, 0)$$

Bien classer tout le monde ($k = 1$) :

$$f(x_i)y_i > 1 \quad i = 1, n$$

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Fidélité et marge



$$f(x_i)y_i > k \quad k > 0, \quad i = 1, n$$

$\mathcal{H} = \{x \mid f(x) = 0\}$: frontière

Marge :

$$\min_{i=1,n} d(\mathcal{H}, x_i) = \min_{i=1,n} \max(1 - f(x_i)y_i, 0)$$

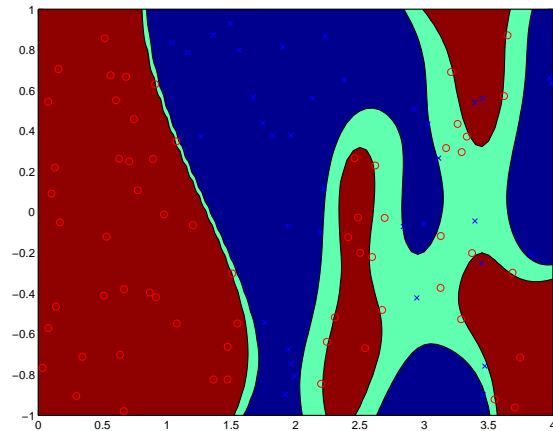
Bien classer tout le monde ($k = 1$) :

$$f(x_i)y_i > 1 \quad i = 1, n$$

1 est la marge minimale

(1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Fidélité et droit à l'erreur : minimiser l'erreur



$$f(x_i)y_i > 1 \quad i = 1, n$$

Introduisons une variable d'écart ξ_i

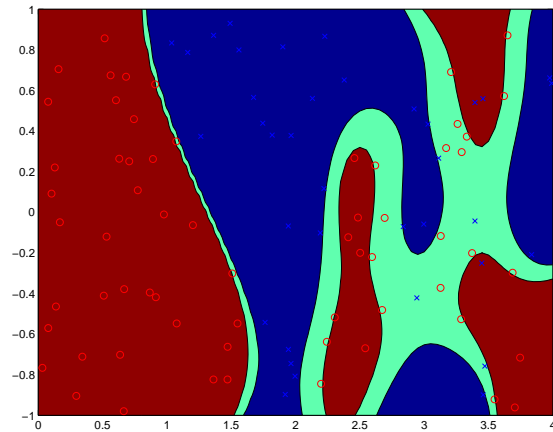
Bien classer **a peu près** tout le monde :

$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

où ξ_i est une variable d'écart

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Fidélité et droit à l'erreur : minimiser l'erreur



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

Introduisons une variable d'écart ξ_i

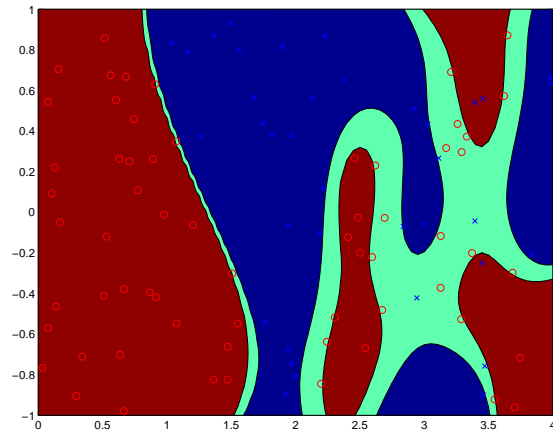
Bien classer **a peu près** tout le monde :

$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

où ξ_i est une variable d'écart

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Fidélité et droit à l'erreur : minimiser l'erreur



$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, i = 1, n$
Introduisons une variable d'écart ξ_i

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

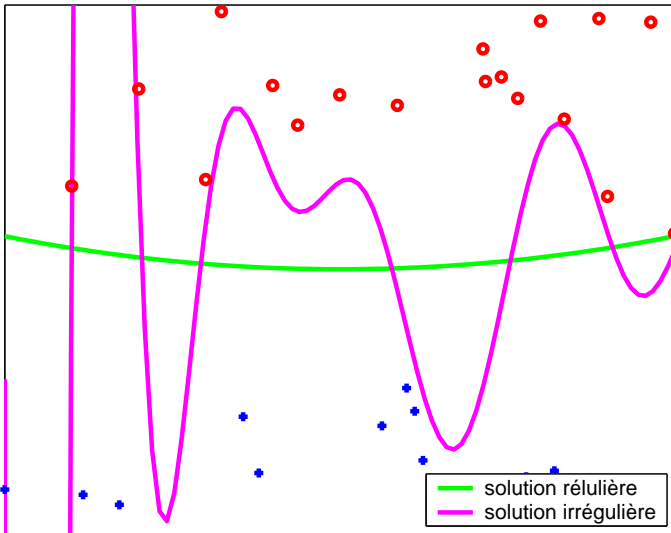
Bien classer **a peu près** tout le monde :

$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, i = 1, n$$

où ξ_i est une variable d'écart $\xi_i = 0$; $\underbrace{\xi_i \geq 1}_{\text{mal classé}} ; 0 < \xi_i < 1$

- (1) **Fidélité** - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Régularité



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

les deux solutions vérifient $\xi_i = 0, \quad i = 1, n$

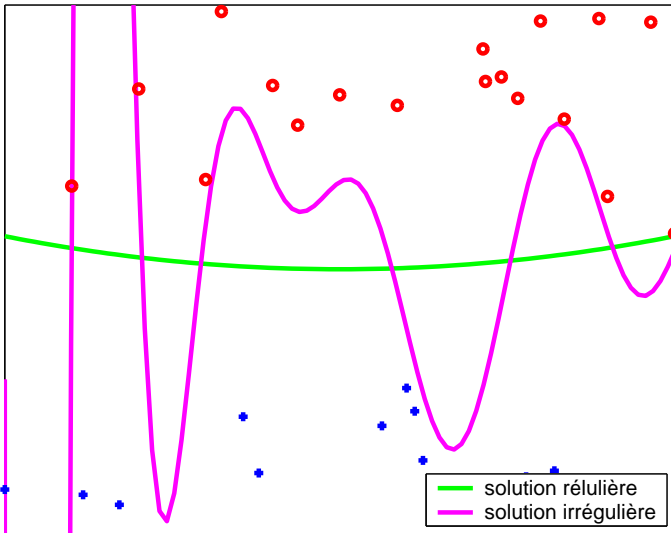
par exemple

■ « l'énergie » de f : la norme de sa dérivée (cf les splines)

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Régularité



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

les deux solutions vérifient $\xi_i = 0, \quad i = 1, n$

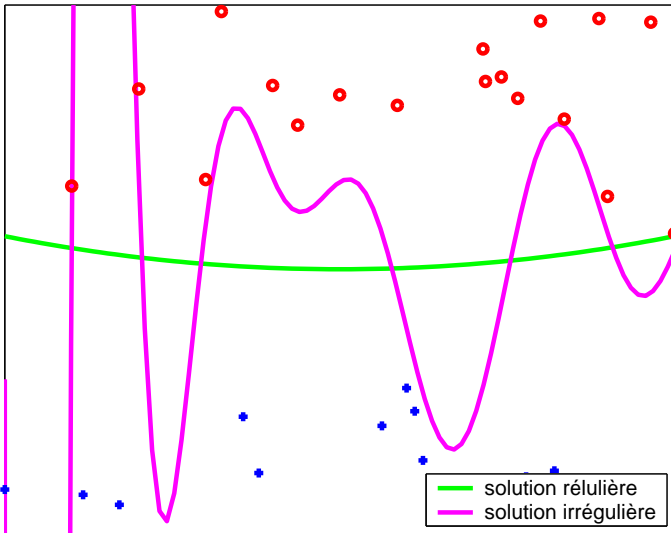
par exemple

- « l'énergie » de f : la norme de sa dérivée (cf les splines)
- la longueur de f - la taille du code calculant f

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Régularité



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

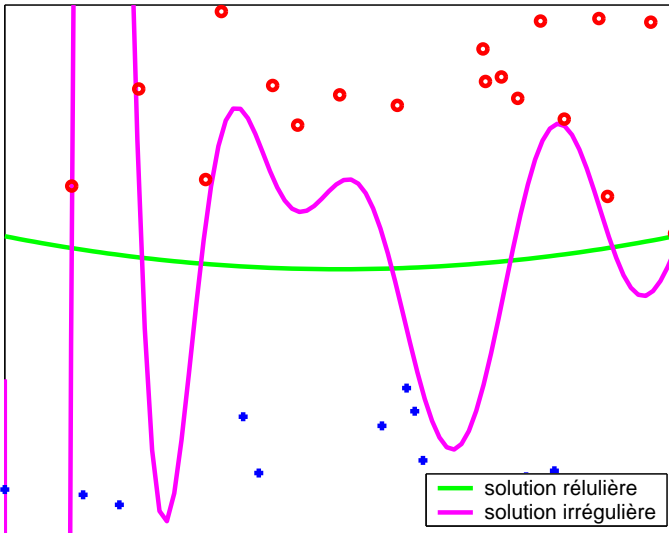
les deux solutions vérifient $\xi_i = 0, \quad i = 1, n$

par exemple

- « l'énergie » de f : la norme de sa dérivée (cf les splines)
- la longueur de f - la taille du code calculant f
- une norme de f au sens de \mathcal{H} (défini a priori) : $\|f\|_{\mathcal{H}}$

(1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Régularité



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

$$\min_f \|f\|_{\mathcal{H}}$$

$$(f(x) = \sum_{j \in J} w_j \phi_j(x) + b)$$

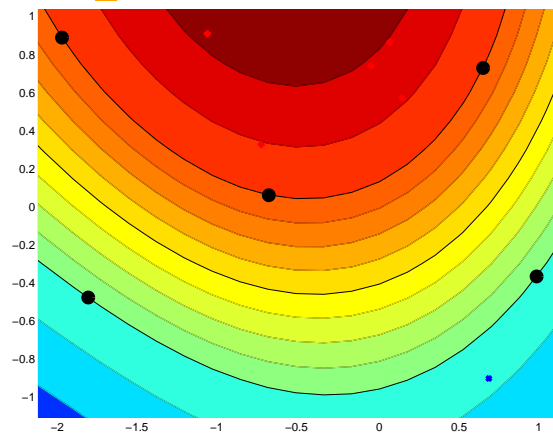
par exemple

- « l'énergie » de f : la norme de sa dérivée (cf les splines)
- la longueur de f - la taille du code calculant f
- une norme de f au sens de \mathcal{H} (défini a priori) : $\|f\|_{\mathcal{H}}$
- un terme de régularisation : une fonctionnelle positive assurant l'unicité de la solution

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Régularité et marge



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

$$\min_f \|f\|_{\mathcal{H}}^2$$

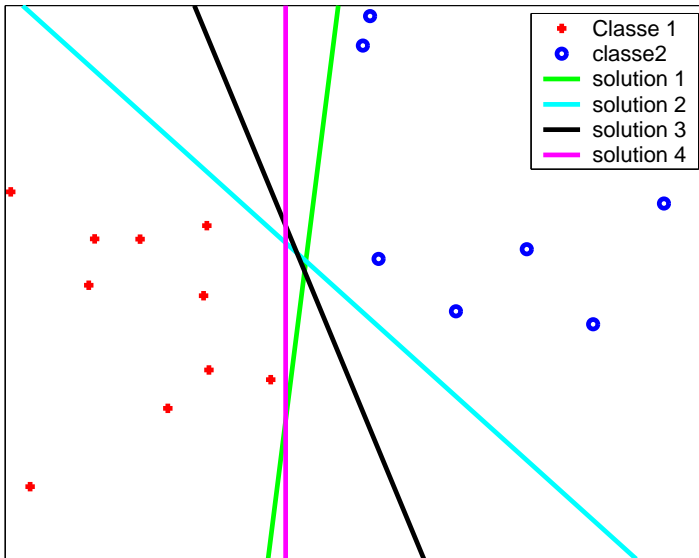
$$\blacksquare \mathbb{P}(err) \leq \underbrace{\sum_{i=1}^n \mathbb{I}_{\{\text{signe}(f(x_i)) \neq y_i\}}}_{\text{Fidélité}} + \varphi\left(\frac{1}{\text{marge}}\right)$$

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Régularité et marge

Cas linéaire : quelle solution choisir ?



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

$$\min_f \|f\|_{\mathcal{H}}^2$$

$$\blacksquare \mathbb{P}(err) \leq \underbrace{\sum_{i=1}^n \mathbb{I}_{\{\text{signe}(f(x_i)) \neq y_i\}}}_{\text{Fidélité}} + \varphi\left(\frac{1}{\text{marge}}\right)$$

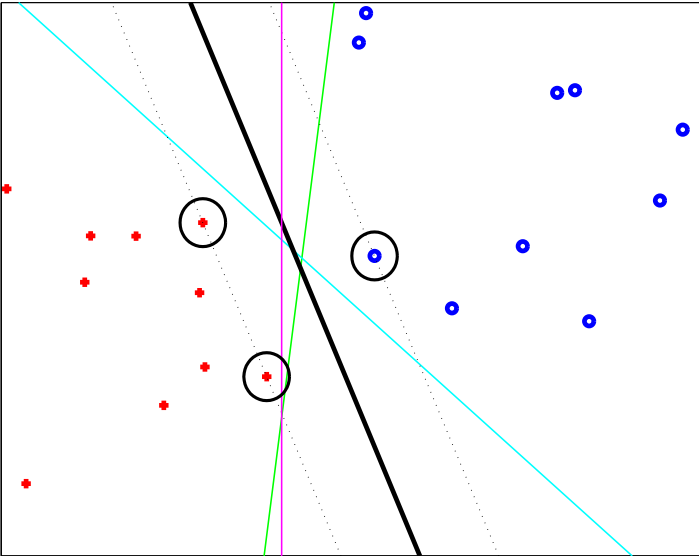
■ minimiser $\mathbb{P}(err) \Leftrightarrow$ maximiser la marge

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Régularité et marge

Celle qui maximise la marge



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

$$\min_f \|f\|_{\mathcal{H}}^2$$

$$\blacksquare \mathbb{P}(err) \leq \underbrace{\sum_{i=1}^n \mathbb{I}_{\{\text{signe}(f(x_i)) \neq y_i\}}}_{\text{Fidélité}} + \varphi\left(\frac{1}{\text{marge}}\right)$$

■ minimiser $\mathbb{P}(err) \Leftrightarrow$ maximiser la marge

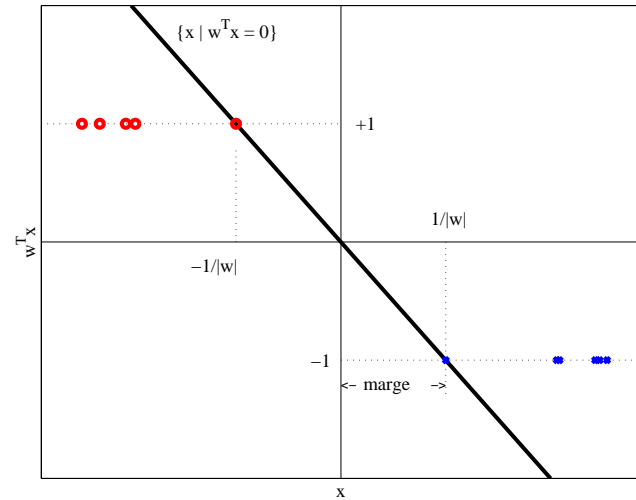
■ maximiser la robustesse \Leftrightarrow maximiser la marge

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Régularité et marge

Valeur de la marge dans le cas monodimensionnel



$$f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n$$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$

$$\min_f \|f\|_{\mathcal{H}}^2 \Leftrightarrow \min_{\mathbf{w}} \sum_{j=1}^{\infty} w_j^2$$

$$\blacksquare \mathbb{P}(err) \leq \underbrace{\sum_{i=1}^n \mathbb{I}_{\{\text{signe}(f(x_i)) \neq y_i\}}}_{\text{Fidélité}} + \varphi\left(\frac{1}{\text{marge}}\right)$$

■ minimiser $\mathbb{P}(err) \Leftrightarrow$ maximiser la marge

■ maximiser la robustesse \Leftrightarrow maximiser la marge

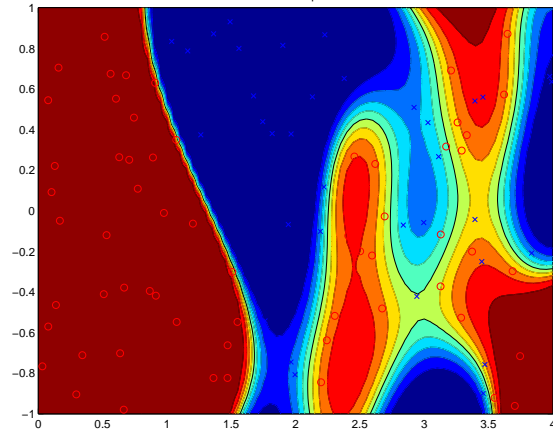
■ maximiser la marge \Leftrightarrow minimiser $\|f\|_{\mathcal{H}}^2$

(1) Fidélité - (3) Décision « locale »

(2) Régularité - (4) Points « frontière »

Apprendre : choisir une hypothèse

Discrimination par SVM



- se donner un ensemble \mathcal{H} assez grand
- trouver $f \in \mathcal{H} : f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, i = 1, n$

$$\min_{\xi_i} \sum_{i=1}^n \xi_i$$
$$\min_f \|f\|_{\mathcal{H}}^2$$

soit $(\phi_j)_{j \in J}$ une base orthonormée de fonctions (polynômes, fourier, ondelettes...)

$$f(x) = \sum_{j=1}^{\infty} w_j \phi_j(x) + b$$

l'ensemble des hypothèses est alors de la forme

$$\mathcal{H} = \left\{ f \mid f(x) = \sum_{j=1}^{\infty} w_j \phi_j(x) + b \right\}$$

f est « linéaire » en ϕ et non linéaire en x

Principe : « qui se ressemble s'assemble »

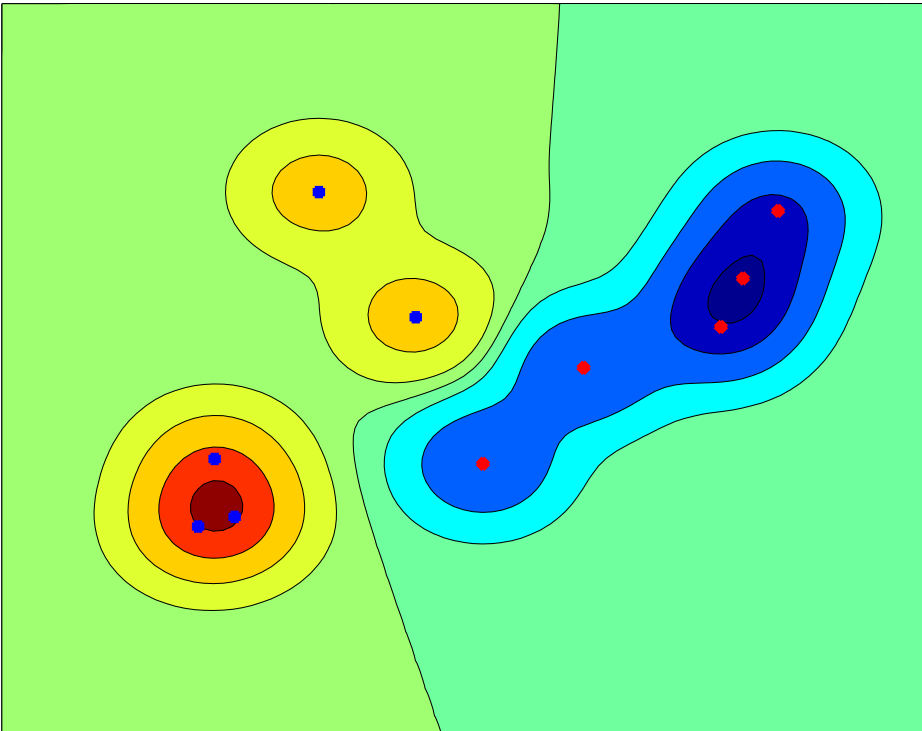
Principes :

- Mesure de similarité
(pas nécessairement symétrique)
- Zone d'influence

Les noyaux :

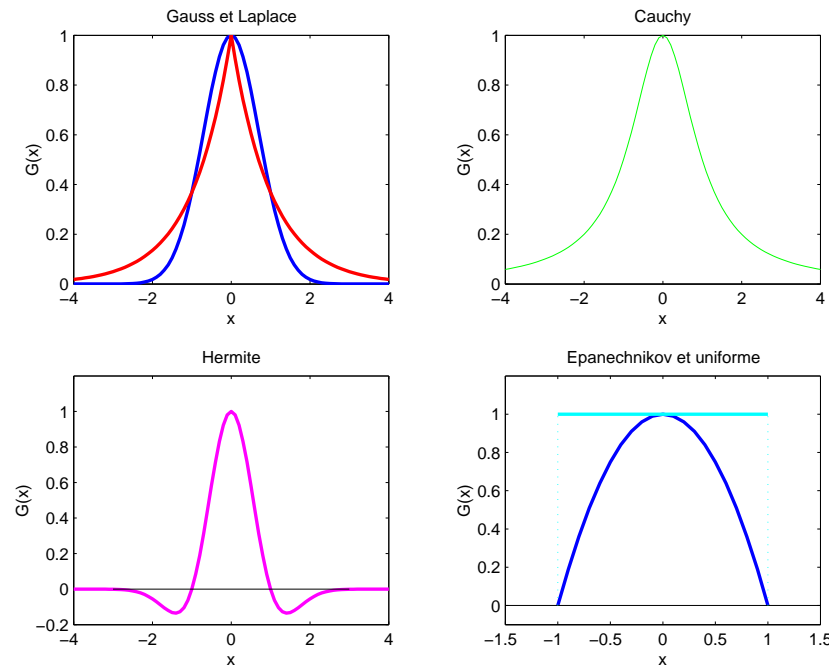
- fonction de deux variables

$$K(x, y)$$



Estimateur à base de **noyaux**

Malédiction de la dimensionnalité



noyaux définis positifs

noyau multidimensionnel produit :

$$K_b(\mathbf{u}, \mathbf{v}) = \prod_{\ell=1}^L K_b(u_\ell, v_\ell)$$

noyaux « radiaux » $\rho = \|\mathbf{u} - \mathbf{v}\|^2$ (distance entre les deux variables)

noyaux « projectifs » $\mathbf{u}^\top \mathbf{v} = \sum_{\ell=1}^L u_\ell v_\ell$

formule de « passage »

$$\|\mathbf{u} - \mathbf{v}\|^2 = \|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 - 2\mathbf{u}^\top \mathbf{v}$$

Quelques exemples de noyaux

le noyau gaussien

$$K_b(u, v) = \frac{1}{Z} \exp^{-\frac{\rho}{b}}$$

le noyau de Cauchy

$$K_b(u, v) = \frac{1}{Z} \frac{1}{1 + \frac{\rho}{b}}$$

le noyau uniforme

$$K_b(u, v) = \frac{1}{Z} \mathbb{I}_{\{\rho \leq b\}}$$

le noyau de Fourier régularisé

$$K_b(u, v) = \frac{1}{Z} \cosh \left(\pi - \frac{|u - v|}{b} \right)$$

le noyau scalaire

$$K_b(\mathbf{u}, \mathbf{v}) = (\mathbf{u}^\top \mathbf{v} + 1)^b$$

le noyau de Laplace

$$K_b(u, v) = \frac{1}{Z} \exp^{-\frac{|u - v|}{b}}$$

le noyau d'Hermite

$$K_b(u, v) = \frac{1}{Z} (b - \rho) \exp^{-\frac{\rho}{b}}$$

le noyau d'Epanechnikov

$$K_b(u, v) = \frac{1}{Z} (b - \rho) \mathbb{I}_{\{\rho \leq b\}}$$

le noyau sigmoïde

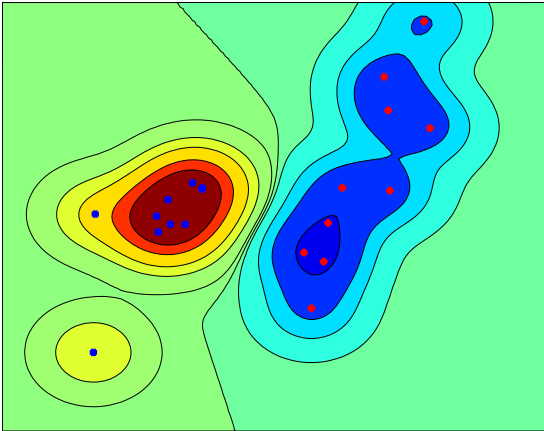
$$K_b(\mathbf{u}, \mathbf{v}) = \frac{1}{Z} \tanh (b(\mathbf{u}^\top \mathbf{v}) + b_0)$$

le noyau de Hardy

$$K_b(u, v) = \frac{1}{(\mathbf{u}^\top \mathbf{v} + 1)^b}$$

noyaux de chaines de caractères, de graphes, d'automates...

Comment choisir \mathcal{H} ?



$$f(x) = \sum_{j=1}^{\infty} w_j \phi_j(x)$$

Influence locale \Rightarrow Noyaux. Influence fixe (Parzen + MAP)

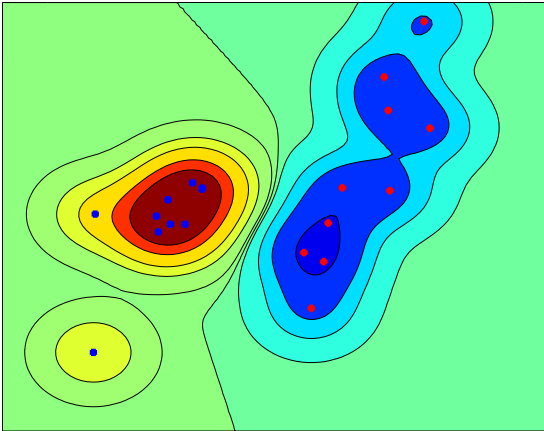
$$f(x) = \sum_{i=1}^n y_i K_b(x, x_i)$$

Influence ajustée

$$f(x) = \sum_{i=1}^n a_i K_b(x, x_i)$$

- (1) Fidélité - (3) **Décision « locale »**
(2) Régularité - (4) Points « frontière »

Comment choisir \mathcal{H} ?



$$f(x) = \sum_{j=1}^{\infty} w_j \phi_j(x)$$

Influence locale \Rightarrow Noyaux. Influence fixe (Parzen + MAP)

$$f(x) = \sum_{i=1}^n y_i K_b(x, x_i)$$

Influence ajustée

$$f(x) = \sum_{i=1}^n a_i K_b(x, x_i)$$

on construit \mathcal{H} à partir du noyau K

- (1) Fidélité - (3) Décision « locale »
- (2) Régularité - (4) Points « frontière »

Comment choisir construire \mathcal{H} (les hypothèses ?)

■ au commencement était le noyau...

■ noyau : $k(x, y) \quad \forall x, y \in \Omega$

■ ...positif $\sum_{i=1}^n \sum_{j=1}^m \alpha_i \alpha_j k(x_i, x_j) > 0$

■ $\mathcal{H}_0 = \{f \in \mathbb{R}^\Omega \mid f(x) = \sum_{i=1}^n \alpha_i k(x, x_i), n \in \mathbb{N}, \alpha_i \in \mathbb{R}, x_i \in \Omega\}$

■ et un produit scalaire on \mathcal{H}_0

$$\langle f(\cdot), g(\cdot) \rangle_{\mathcal{H}_0} = \sum_{i=1}^n \sum_{j=1}^m \alpha_i \beta_j k(x_i, x_j)$$

■ propriétés du produit scalaire sur \mathcal{H}_0

■ $\langle f(\cdot), k(x, \cdot) \rangle_{\mathcal{H}_0} = f(x)$

■ $\langle k(x, \cdot), k(y, \cdot) \rangle_{\mathcal{H}_0} = k(x, y)$

Evaluation

(reproduction)

le noyau représente une grande partie de la connaissance a priori

L'astuce du noyau

Théorème de Mercer : Si K est un noyau défini positif, il existe une famille $(\phi_j)_{j \in \mathcal{J}}$ orthonormée telle que :

$$K_b(\mathbf{x}, \mathbf{y}) = \sum_{j \in J} \phi_j(\mathbf{x}) \phi_j(\mathbf{y}) \quad (1)$$

toute fonction $f \in \mathcal{H}$ s'écrit alors :

$$f(\mathbf{x}) = \sum_{j \in J} w_j \phi_j(\mathbf{x}) = \sum_{i=1}^n a_i K_b(\mathbf{x}, \mathbf{x}_i)$$

$$\|f\|_{\mathcal{H}}^2 = \mathbf{w}^\top \mathbf{w} = \mathbf{a}^\top K \mathbf{a}$$

où K est la matrice d'influence.

$$K_{ij} = K_b(\mathbf{x}_i, \mathbf{x}_j)$$

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

L'astuce du noyau

Théorème de Mercer : Si K est un noyau défini positif, il existe une famille $(\phi_j)_{j \in \mathcal{J}}$ orthonormée telle que :

$$K_b(\mathbf{x}, \mathbf{y}) = \sum_{j \in J} \phi_j(\mathbf{x}) \phi_j(\mathbf{y}) \quad (2)$$

toute fonction $f \in \mathcal{H}$ s'écrit alors :

$$f(\mathbf{x}) = \sum_{j \in J} w_j \phi_j(\mathbf{x}) = \sum_{i=1}^n a_i K_b(\mathbf{x}, \mathbf{x}_i)$$

$$\|f\|_{\mathcal{H}}^2 = \mathbf{w}^\top \mathbf{w} = \mathbf{a}^\top K \mathbf{a}$$

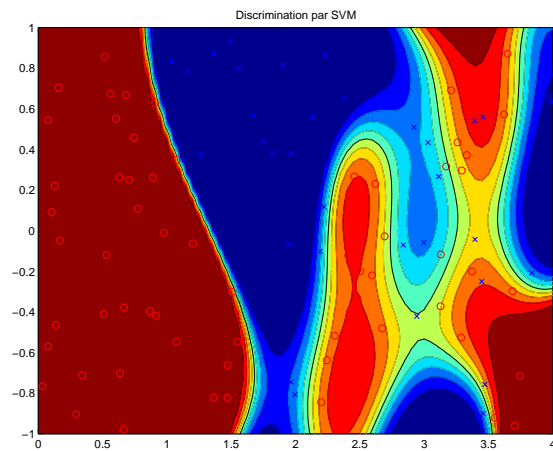
où K est la matrice d'influence.
 $K_{ij} = K_b(\mathbf{x}_i, \mathbf{x}_j)$

dim ∞

dim n

(1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Quel critère minimiser ?



$$\left. \begin{array}{l} f(x_i)y_i > 1 - \xi_i \quad \xi_i > 0, \quad i = 1, n \\ \min_{\xi_i} \sum_{i=1}^n \xi_i \end{array} \right\} \text{Fidélité}$$

$$\min_{\mathbf{w}} \sum_{j=1}^{\infty} w_j^2 \quad \dots \quad \text{Régularité}$$

Comment choisir la fonction de discrimination ?

$$f(x) = \sum_{j=1}^{\infty} w_j \phi_j(x)$$

- (1) Fidélité - (3) Décision « locale »
- (2) Régularité - (4) Points « frontière »

Ensemble d'hypothèses + Critère = Le problème SVM

$$\mathcal{H} = \left\{ f : \mathbb{R}^L \rightarrow \mathbb{R} \mid \exists \mathbf{a}, \mathbf{c} ; f(\mathbf{x}) = \sum_{j=1}^m c_j \varphi_j(\mathbf{x}) + \sum_{\ell=1}^{n_{\text{sup}}} a_{\ell} K_b(\mathbf{x}, \mathbf{x}_{\ell}) \right\}$$

où n_{sup} est le nombre de vecteurs supports
problème de minimisation sous contraintes :

$$\left\{ \begin{array}{ll} \min_{\mathbf{w}} & \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{avec} & y_i f(\mathbf{x}_i) > 1 - \xi_i \quad i = 1, n \\ \text{et} & \xi_i > 0 \quad i = 1, n \end{array} \right. \quad (3)$$

$$\text{où : } f(\mathbf{x}) = \sum_{k=1}^{\infty} w_k \phi_k(\mathbf{x}) + \sum_{j=1}^m c_j \varphi_j(\mathbf{x})$$

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

Ensemble d'hypothèses + Critère = Le problème SVM

$$\mathcal{H} = \left\{ f : \mathbb{R}^L \rightarrow \mathbb{R} \mid \exists \mathbf{a}, \mathbf{c} ; f(\mathbf{x}) = \sum_{j=1}^m c_j \varphi_j(\mathbf{x}) + \sum_{\ell=1}^{n_{\text{sup}}} a_{\ell} K_b(\mathbf{x}, \mathbf{x}_{\ell}) \right\}$$

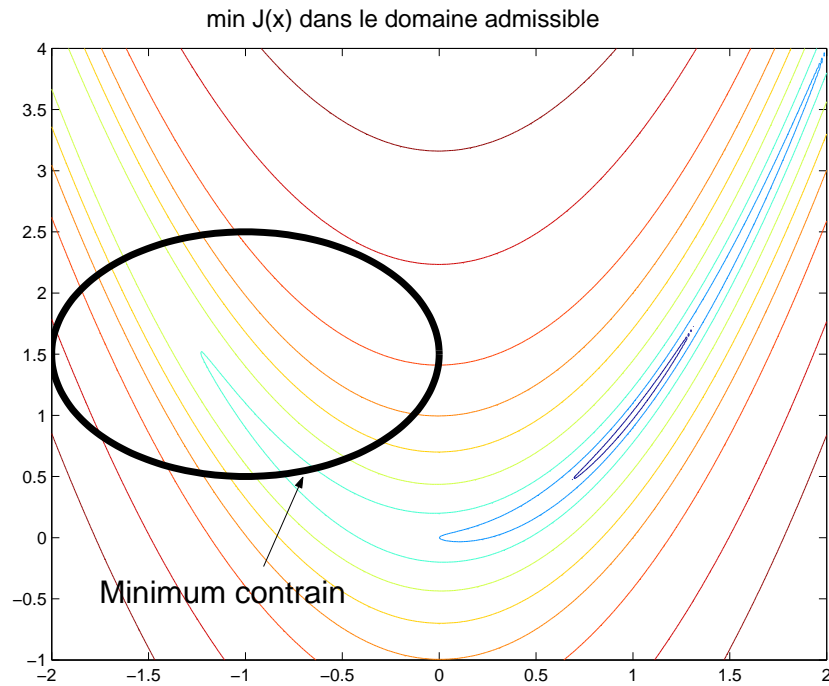
où n_{sup} est le nombre de vecteurs supports
problème de minimisation sous contraintes :

$$\left\{ \begin{array}{l} \min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{avec} \\ \text{et} \quad y_i f(\mathbf{x}_i) > 1 - \xi_i \quad i = 1, n \\ \xi_i > 0 \quad i = 1, n \end{array} \right. \quad (4)$$

où : $f(\mathbf{x}) = \sum_{k=1}^{\infty} w_k \phi_k(\mathbf{x}) + \sum_{j=1}^m c_j \varphi_j(\mathbf{x})$

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - (4) Points « frontière »

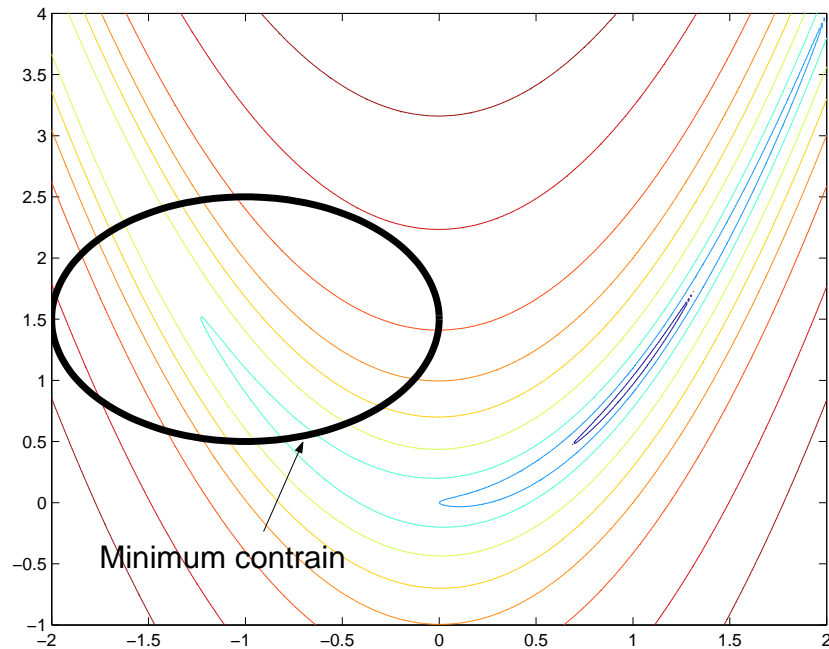
Minimisation sous contraintes (cas séparable)



$$\left\{ \begin{array}{l} \min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{avec} \quad y_i f(\mathbf{x}_i) > 1 \quad i = 1, n \end{array} \right.$$

Minimisation sous contraintes (cas séparable)

min J(x) dans le domaine admissible



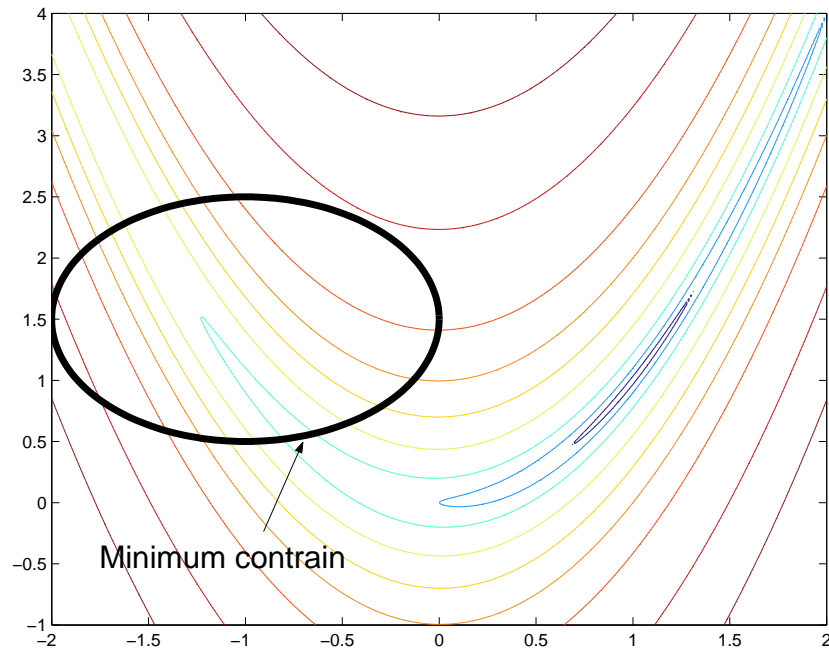
$$\left\{ \begin{array}{l} \min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{avec} \quad y_i f(\mathbf{x}_i) > 1 \quad i = 1, n \end{array} \right.$$
$$\Leftrightarrow$$

$$\min_{\mathbf{w}, \mathbf{c}} \max_{\lambda} \mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda) \quad \text{Lagrangien}$$

$$\mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \underbrace{\sum_{i=1}^n}_{\text{les exemples}} \lambda_i (y_i f(\mathbf{x}_i) - 1)$$

Minimisation sous contraintes (cas séparable)

min $J(\mathbf{x})$ dans le domaine admissible



$$\left\{ \begin{array}{l} \min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} \\ \text{avec} \quad y_i f(\mathbf{x}_i) > 1 \quad i = 1, n \end{array} \right.$$

\Leftrightarrow

$$\min_{\mathbf{w}, \mathbf{c}} \max_{\lambda} \mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda) \quad \text{Lagrangien}$$

$$\mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \underbrace{\sum_{i=1}^n}_{\text{les exemples}} \lambda_i (y_i f(\mathbf{x}_i) - 1)$$

Multiplicateur de Lagrange $\lambda_i =$ **influence de l'exemple i dans la solution**

interprétation : $\lambda_i = 0 \rightarrow$ pas d'influence $\lambda_i > 0 \rightarrow$ exemple *support*

Reformulation dans l'espace des exemples

$$\mathcal{L}(\mathbf{w}, \mathbf{a}, \lambda) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \lambda_i (y_i f(\mathbf{x}_i) - 1)$$

dont on tire les conditions de Kuhn et Tucker :

$$\begin{cases} \frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda)}{\partial \mathbf{w}} = 0 \\ \frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda)}{\partial \mathbf{c}} = 0 \end{cases} \Leftrightarrow \begin{cases} \mathbf{w} - \sum_{i=1}^n \lambda_i y_i \phi(\mathbf{x}_i) = 0 \\ \sum_{i=1}^n \lambda_i y_i \varphi(\mathbf{x}_i) = 0 \end{cases}$$

conséquence pour f :

$$\begin{aligned} f(\mathbf{x}) = \sum_{k=1}^{\infty} w_k \phi_k(\mathbf{x}) &= \sum_{k=1}^{\infty} \left(\sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i) \right) \phi_k(\mathbf{x}) \\ &= \sum_{i=1}^N \underbrace{\lambda_i y_i}_{a_i} \underbrace{\sum_{k=1}^{\infty} \phi_k(\mathbf{x}) \phi(\mathbf{x}_i)}_{K_b(\mathbf{x}, \mathbf{x}_i)} \end{aligned}$$

Reformulation dans l'espace des exemples

$$\mathcal{L}(\mathbf{w}, \mathbf{a}, \lambda) = \frac{1}{2} \|f\|^2 - \sum_{i=1}^n \lambda_i (y_i f(\mathbf{x}_i) - 1)$$

dont on tire les conditions de Kuhn et Tucker :

$$\begin{cases} \frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda)}{\partial \mathbf{w}} = 0 \\ \frac{\partial \mathcal{L}(\mathbf{w}, \mathbf{c}, \lambda)}{\partial \mathbf{c}} = 0 \end{cases} \Leftrightarrow \begin{cases} \mathbf{w} - \sum_{i=1}^n \lambda_i y_i \phi(\mathbf{x}_i) = 0 \\ \sum_{i=1}^n \lambda_i y_i \varphi(\mathbf{x}_i) = 0 \end{cases}$$

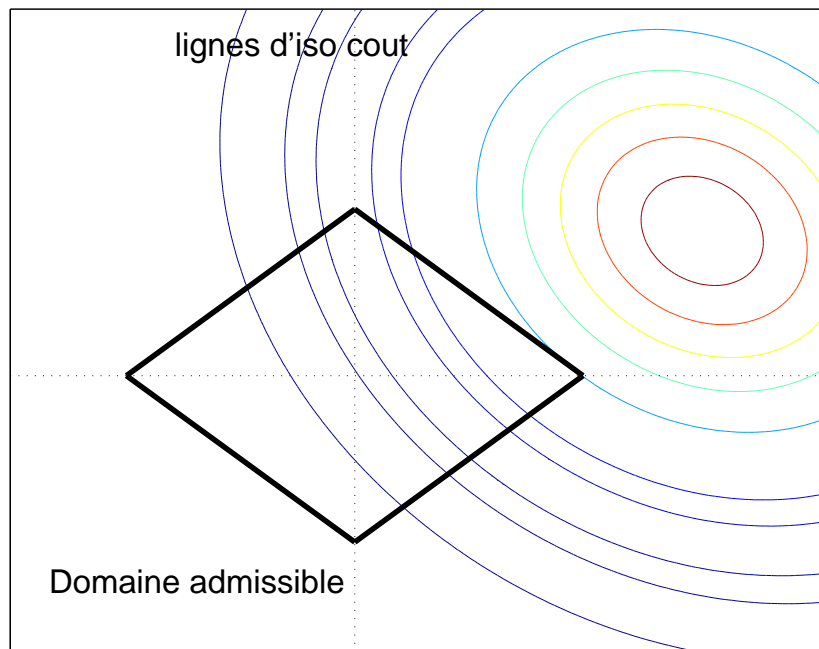
conséquence pour f :

$$\begin{aligned} f(\mathbf{x}) &= \sum_{k=1}^{\infty} w_k \phi_k(\mathbf{x}) = \sum_{k=1}^{\infty} \left(\sum_{i=1}^N \lambda_i y_i \phi(\mathbf{x}_i) \right) \phi_k(\mathbf{x}) \\ &= \sum_{i=1}^N \underbrace{\lambda_i y_i}_{a_i} \underbrace{\sum_{k=1}^{\infty} \phi_k(\mathbf{x}) \phi(\mathbf{x}_i)}_{K_b(\mathbf{x}, \mathbf{x}_i)} \end{aligned}$$

Stratégie

calcul de K
calcul des λ
calcul des a
calcul de f

calcul des λ : problème Dual (2)



$$\left\{ \begin{array}{l} \min_{\lambda} \quad \frac{1}{2} \lambda^{\top} H \lambda + \mathbf{c}^{\top} \lambda \\ \text{avec} \quad \sum_{i=1}^N \lambda_i y_i \varphi_j(x_i) = 0 \quad j = 1, m \\ \text{et} \quad 0 \leq \lambda_i \leq C \quad i = 1, n \end{array} \right.$$

où H est la matrice de terme général $H_{ij} = y_i y_j K_b(\mathbf{x}_i, \mathbf{x}_j)$
et \mathbf{c} un vecteur de 1.

Reformulation de Girosi (97)

$$\min_{\mathbf{a}} \|f(\mathbf{x}_i) - y_i\|_{\mathcal{H}}^2 + \mu \sum_{i=1}^n |a_i|$$

- (1) Fidélité - (3) Décision « locale »
(2) Régularité - **(4) Points « frontière »**

Solution pratique : Problème d'optimisation

- dans le pire des cas...Simplex
- Solution « hors lignes » :
 - Contraintes actives ($\mathcal{O}(n^{1.5})$)
 - asi.insa-rouen.fr/~gloosli
 - asi.insa-rouen.fr/~arakotom
 - Points intérieurs : lent
 - stochastiques : SMO
 - libsvm rapide et complet
 - coresvm (pas très fiable)
- Solution « en ligne » : LaSVM
 - La SVM : très très rapides - $8 \cdot 10^6$ exemples...

si n est grand : plus efficace que les autres méthodes

Conclusion

- Les méthodes à noyaux
 - approximateur universel
 - minimum global unique
 - Parcimonieux (des coef. = 0)

⇒ très rapide (en général)
- Classification : les SVM
 - SVM vs Réseaux de neurones (PMC) : optimisation
 - SVM vs Parzen : parcimonie et vitesse (L^2 vs L^1)
 - SVM vs régression logistique : parcimonie et vitesse
 - SVM : des résultats
- Régression : le kLAR
 - kLAR vs Réseaux de neurones (PMC) : optimisation
 - kLAR vs Noyaux (FBR) : optimisation
 - kLAR vs splines : parcimonie et vitesse (L^2 vs L^1)
- autres applications : one class SVM, kACP, kPLS...

Références

■ SVM

- V. Vapnik : The Nature of Statistical Learning Theory. Springer, 1995.
- N. Cristianini and J. Shawe-Taylor : An Introduction to Support Vector Machines. Cambridge University Press, Cambridge, UK, 2000
- B. Scholkopf et A. Smola : Kernel Machines, 2002

■ Reconnaissance des formes statistiques

- R. O. Duda, P. E. Hart and D. G. Stork : Pattern Classification (2nd ed.), John Wiley and Sons, 2001.
- T. Hastie, R. Tibshirani, and J. Friedman : The Elements of Statistical Learning : Data Mining, Inference, and Prediction, Springer-Verlag, 2001

■ et sur le réseau

- <http://kernel-machines.org>
- <http://www.ph.tn.tudelft.nl/PRInfo/>
- <http://citeseer.nj.nec.com/>
- <http://asi.insa-rouen.fr/~scanu>