

Apprentissage, reproduction, régularisation et noyaux séparateur à vaste marge (SVM)

Stéphane Canu, Xavier Mary et Alain Rakotomamonjy

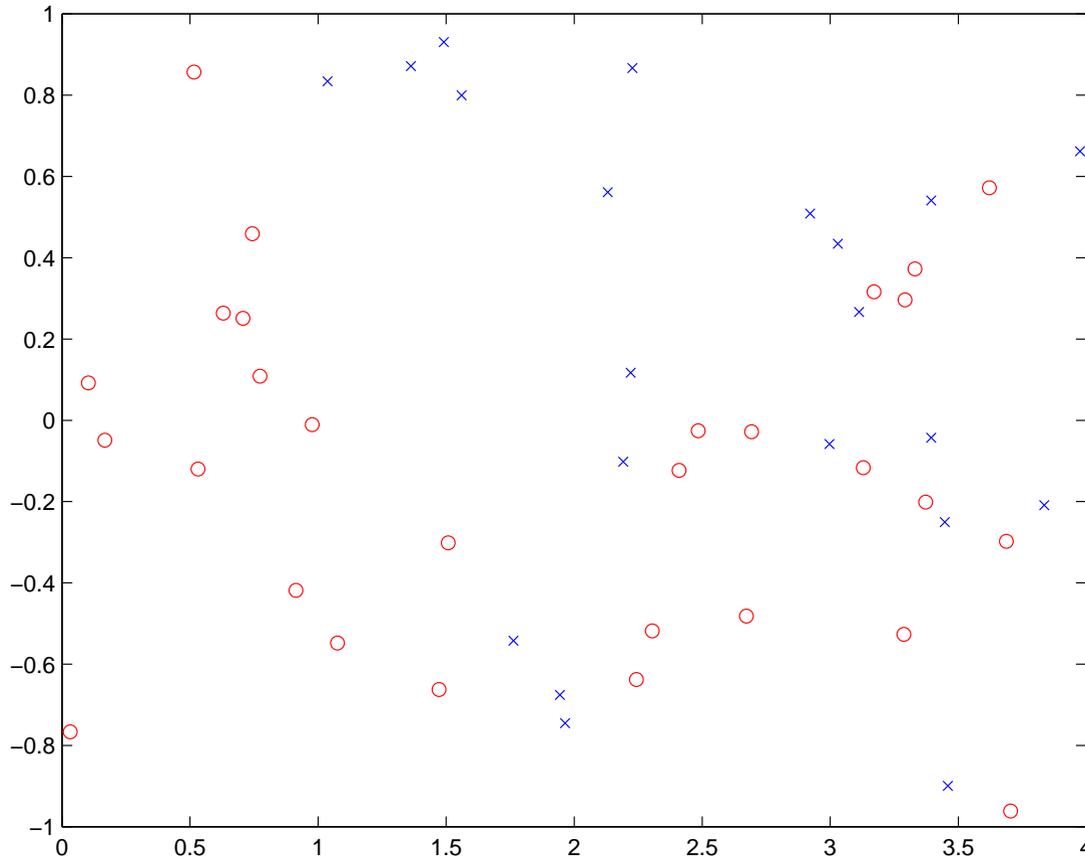
`stephane.canu@insa-rouen.fr`

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

INSA Rouen - Département ASI

Laboratoire PSI, FRE CNRS 2645

A la recherche d'une règle de décision universelle



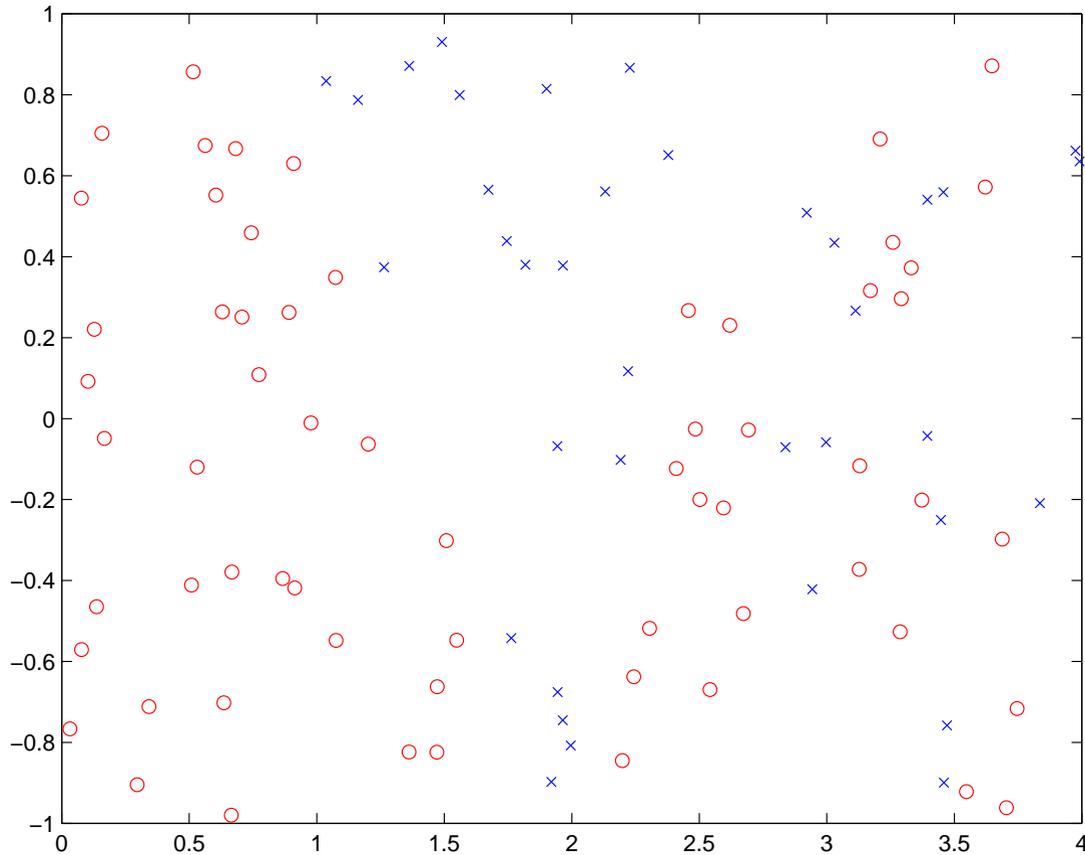
on cherche un algorithme \mathcal{A}
capable de résoudre
tous les problèmes

l'échantillon $(x_i, y_i)_{i=1, n}$

$$\underbrace{\mathbb{P}(err(\mathcal{A}, x_i, y_i))}_{\text{erreur de } \mathcal{A}} \xrightarrow{n \rightarrow \infty} \underbrace{\mathbb{P}_b(err)}_{\text{erreur de bayes}}$$

Tracez la frontière de décision entre ces deux classes ?

Introduction

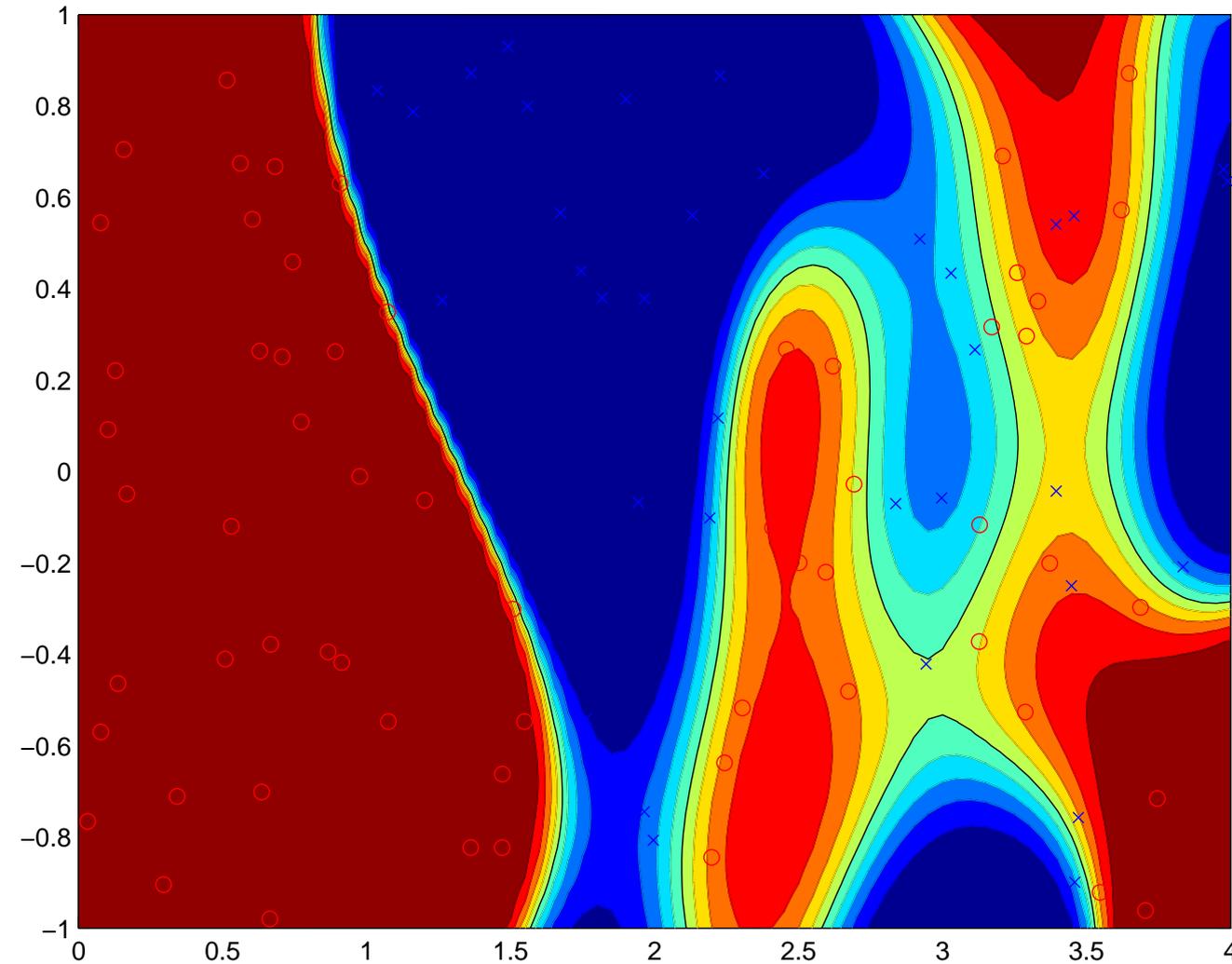


c'est plus facile...
...avec un peu plus de points

Tracez la frontière de décision entre ces deux classes ?

Introduction : le plan

Discrimination par SVM

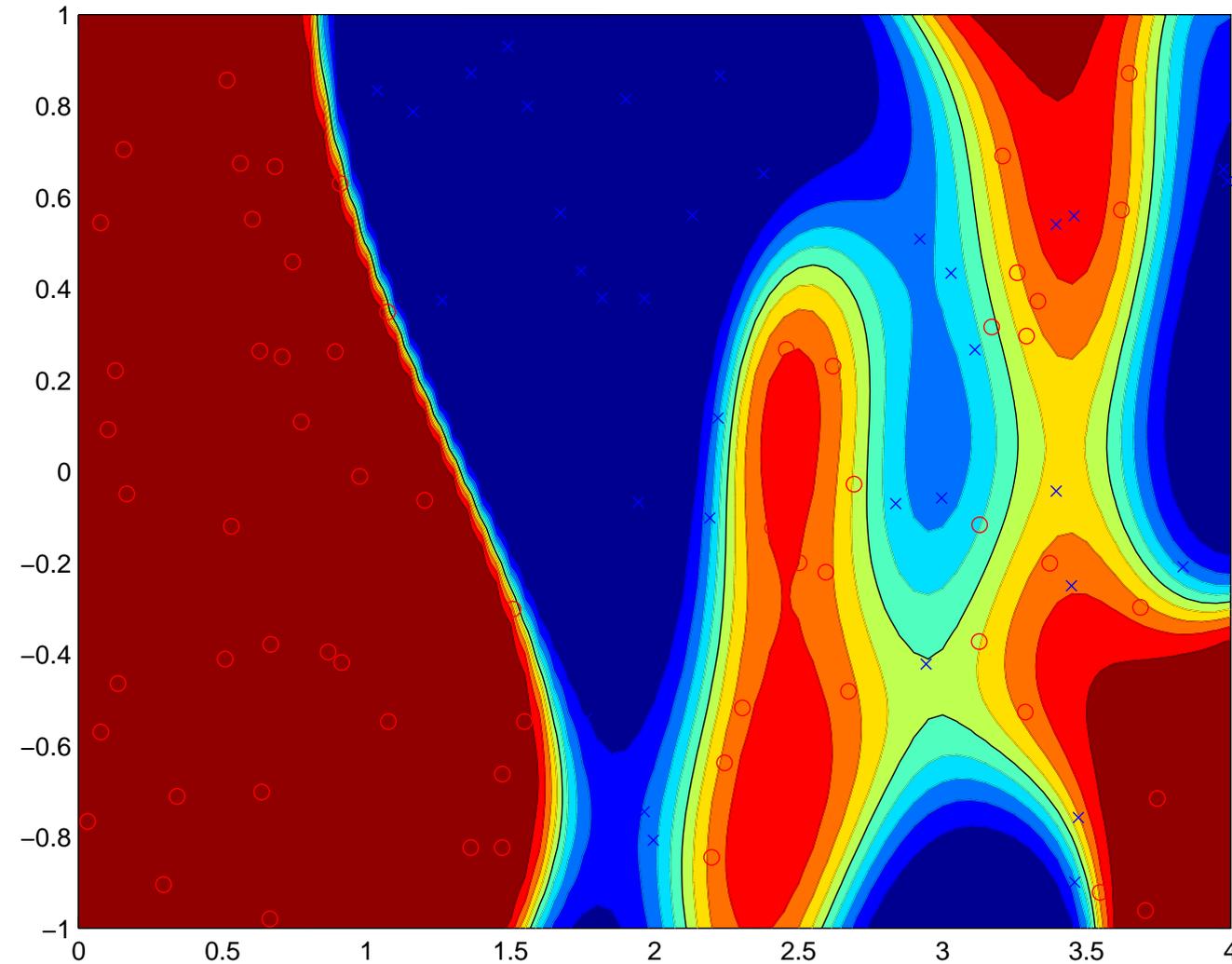


Une solution
⇒ Quels critères ?

■ (1) Fidélité

Introduction : le plan

Discrimination par SVM

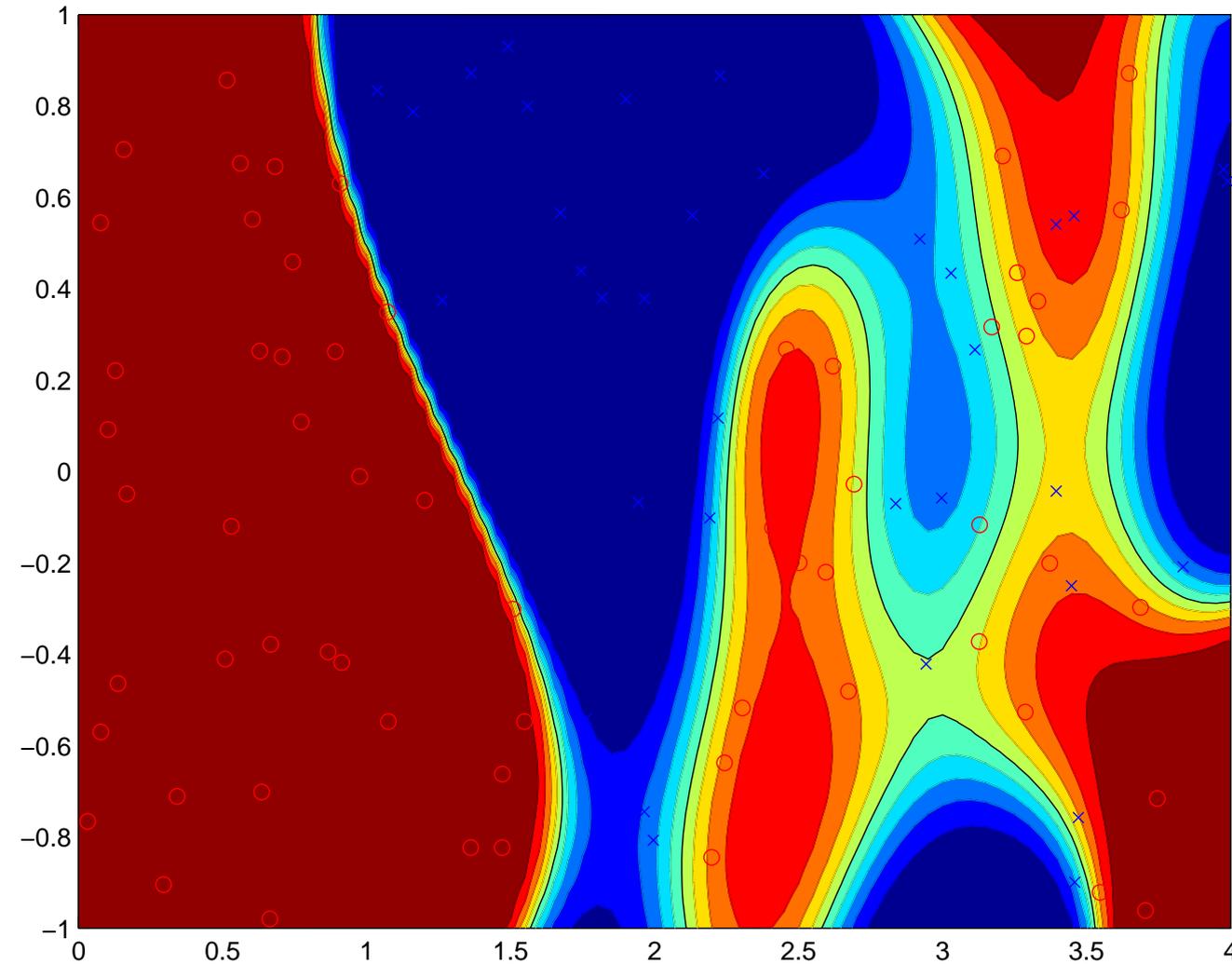


Une solution
⇒ Quels critères ?

- (1) Fidélité
- (2) Régularité

Introduction : le plan

Discrimination par SVM

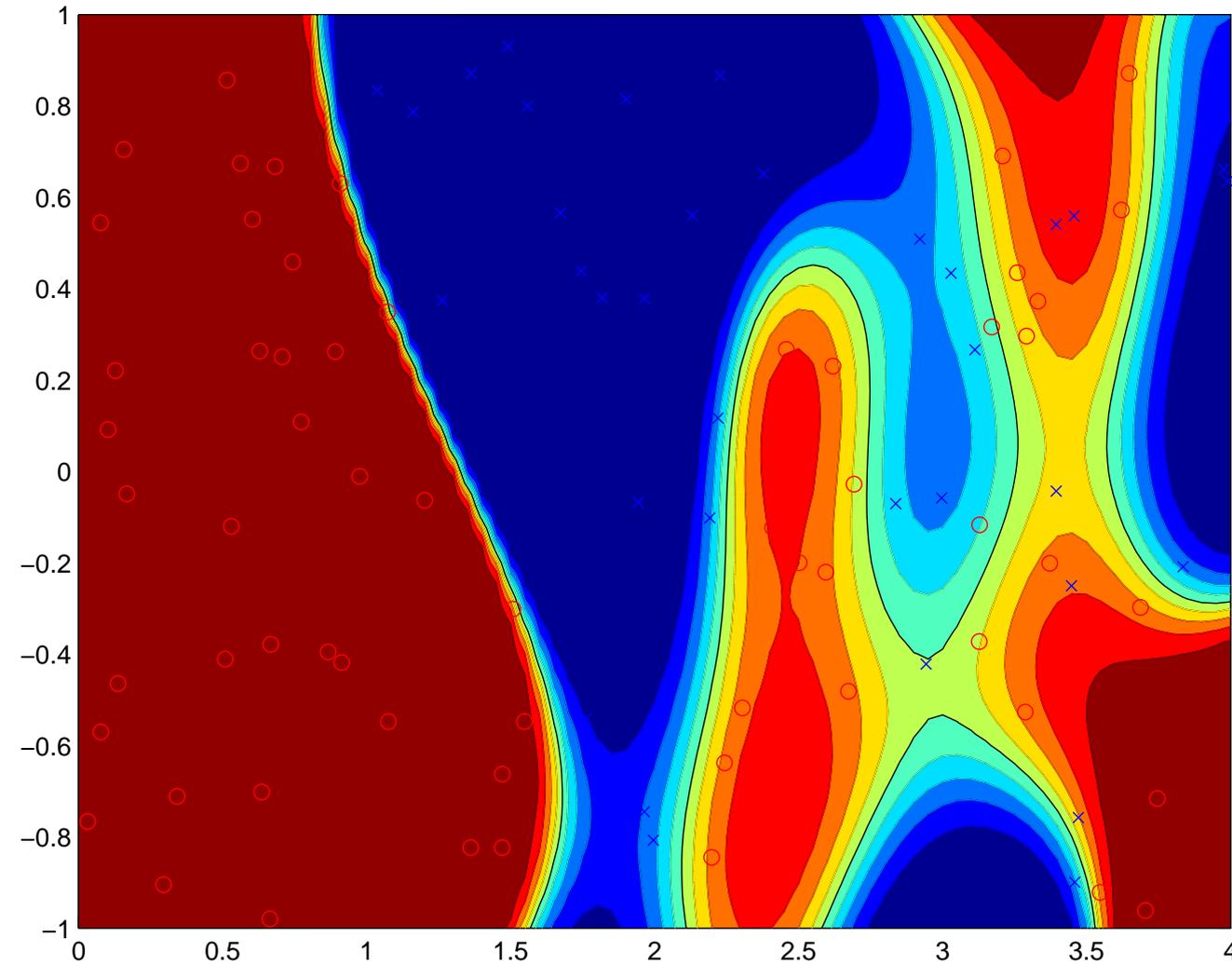


Une solution
⇒ Quels critères ?

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

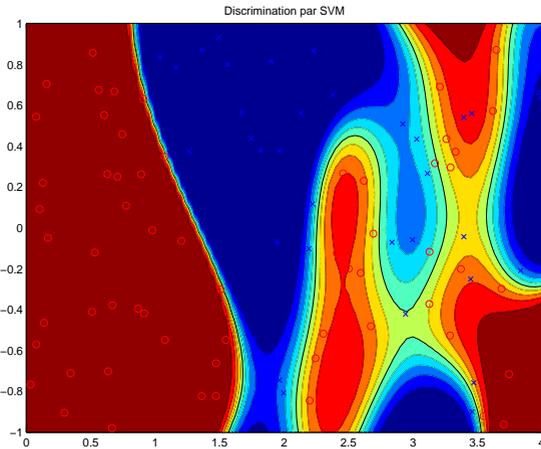
Introduction : le plan

Discrimination par SVM



Une solution
⇒ Quels critères ?

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



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)

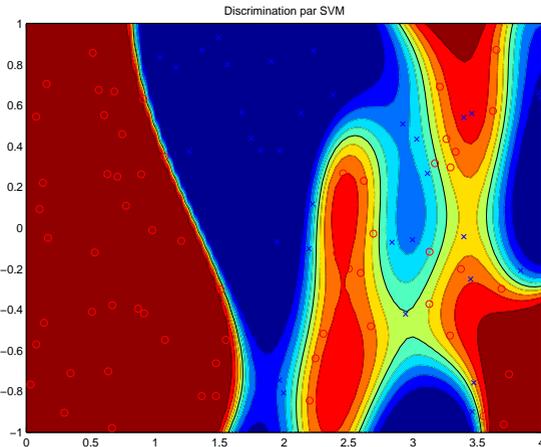
$\mathcal{F} = \{x \mid f(x) = 0\}$: frontière de décision.

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 »

Fidélité



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)

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

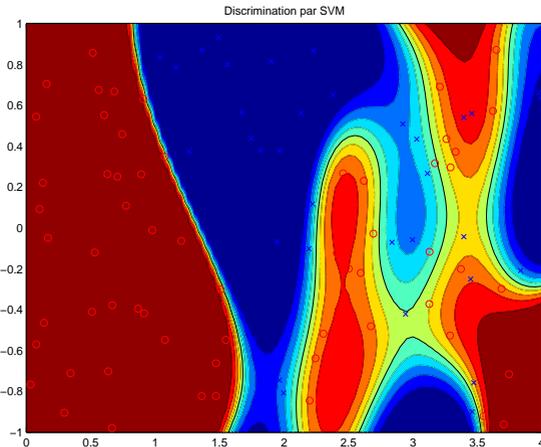
Bien classer tout le monde :

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

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

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

Fidélité



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)

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

Bien classer tout le monde :

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

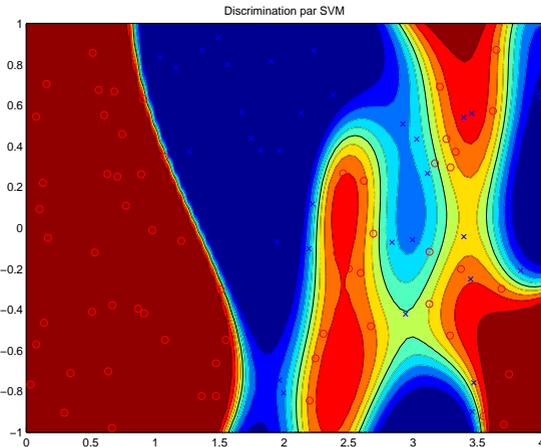
$f(x_i)y_i \geq 0 \quad i = 1, n \quad \text{solution triviale } f = 0$

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

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

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

Fidélité



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)

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

Bien classer tout le monde :

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

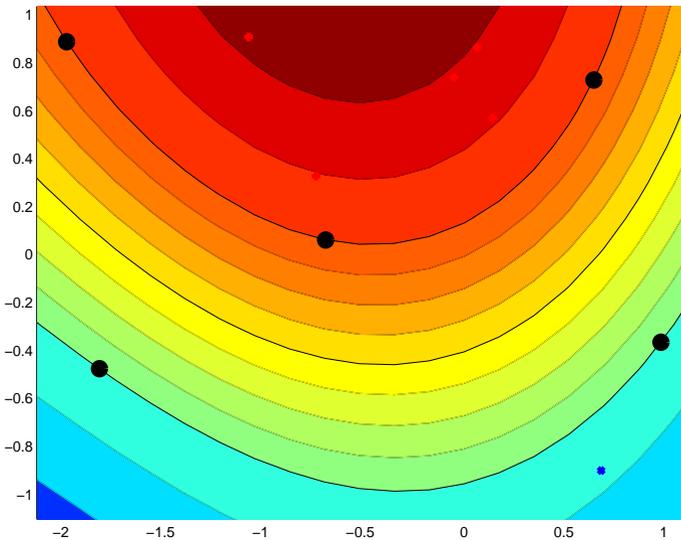
$f(x_i)y_i \geq 0$ $i = 1, n$ solution triviale $f = 0$

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

Marge

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

Fidélité et marge



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

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

Marge :

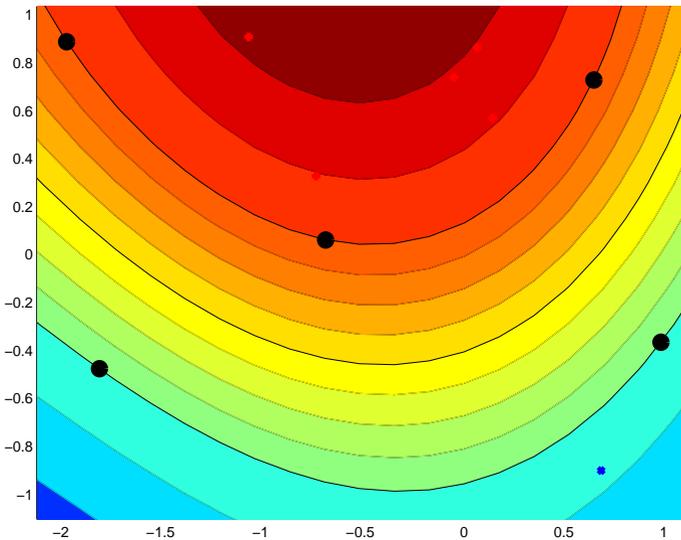
$$\min_{i=1,n} d(\mathcal{F}, x_i) = \min_{i=1,n} f(x_i)y_i$$

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



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

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

Marge :

$$\min_{i=1,n} d(\mathcal{F}, x_i) = \min_{i=1,n} f(x_i)y_i$$

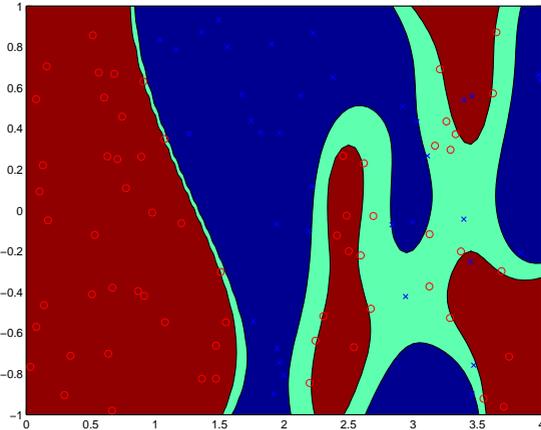
Bien classer tout le monde ($k = 1$) :

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

1 est la marge minimale

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

Fidélité et droit à l'erreur



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

Introduisons une variable d'écart ξ_i

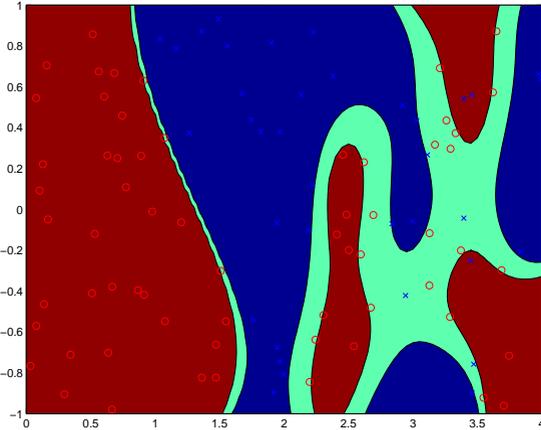
Bien classer **a peu près** tout le monde :

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

où ξ_i est une variable d'écart

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

Fidélité et droit à l'erreur



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

Introduisons une variable d'écart ξ_i

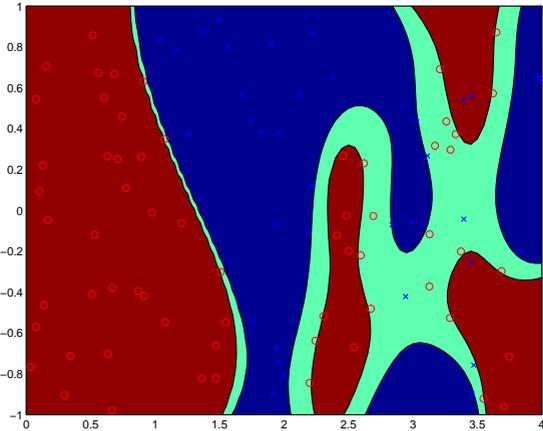
Bien classer **a peu près** tout le monde :

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

où ξ_i est une variable d'écart

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

Fidélité et droit à l'erreur



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

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

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

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

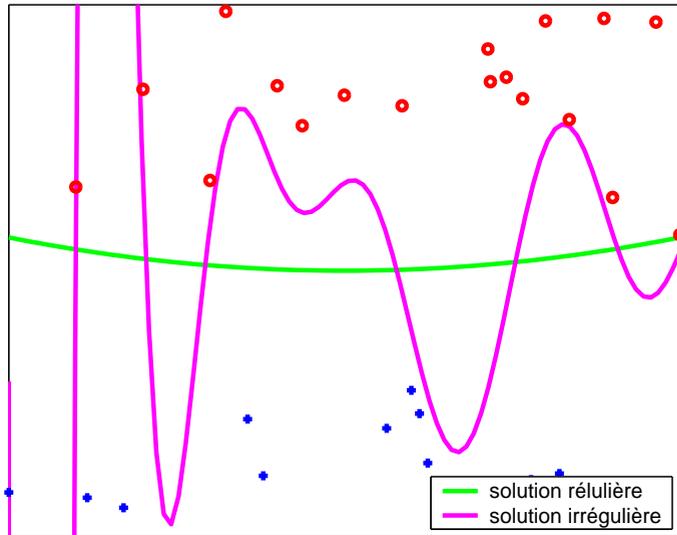
où ξ_i est une variable d'écart $\xi_i = 0;$ $\underbrace{\xi_i \geq 1}$; $0 < \xi_i < 1$

mal classé

(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 $\xi_i = 0, \quad i = 1, n$

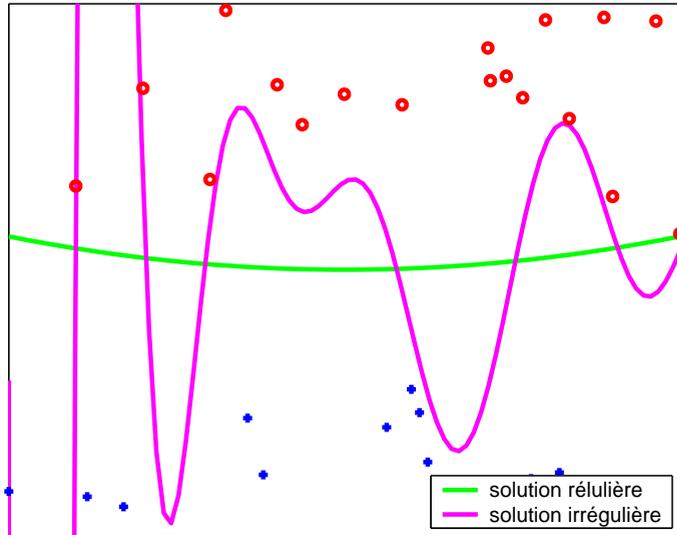
par exemple

- « 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 $\xi_i = 0, \quad i = 1, n$

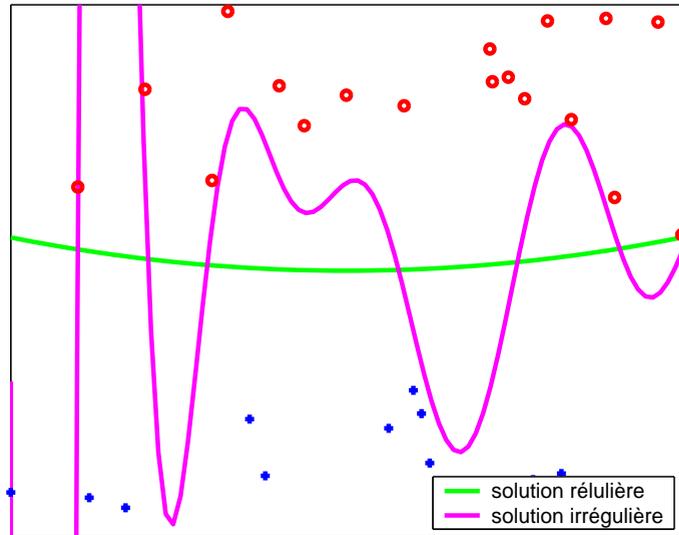
par exemple

- « 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 $\xi_i = 0, \quad i = 1, n$

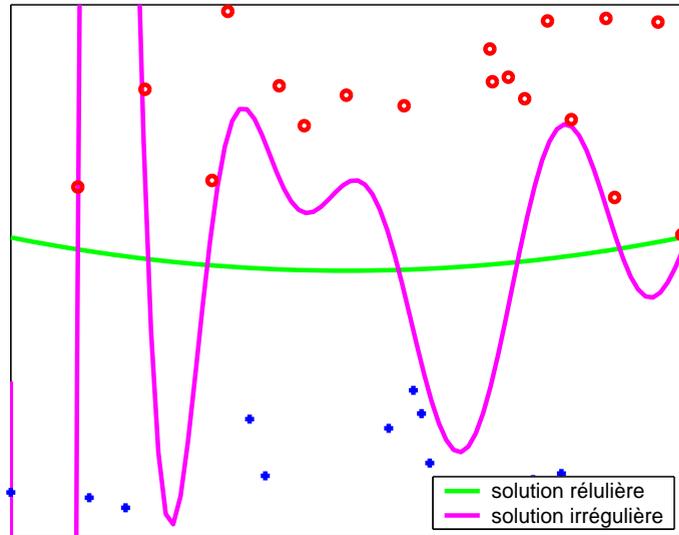
par exemple

- « 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



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

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

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

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 »

Comment choisir \mathcal{H} (l'ensemble des hypothèses ?)

1. \mathcal{H} doit être un espace vectoriel

Comment choisir \mathcal{H} (l'ensemble des hypothèses ?)

1. \mathcal{H} doit être un espace vectoriel
2. $f(x)$ doit avoir un sens...

$$\begin{aligned} \delta_x : \mathcal{H} &\longrightarrow \mathbb{R} \\ f &\longmapsto \delta_x f = f(x) \end{aligned}$$

L^2 ne convient pas !

Comment choisir \mathcal{H} (l'ensemble des hypothèses ?)

1. \mathcal{H} doit être un espace vectoriel
2. $f(x)$ doit avoir un sens...

$$\begin{aligned}\delta_x : \mathcal{H} &\longrightarrow \mathbb{R} \\ f &\longmapsto \delta_x f = f(x)\end{aligned}$$

L^2 ne convient pas !

3. la fonctionnelle d'évaluation doit être continue ($\delta_x \in \mathcal{H}'$) :

$$f_n \xrightarrow{\mathcal{H}} t \implies \forall x \in \mathcal{X}, f_n(x) \xrightarrow{\mathbb{R}} t(x)$$

dans le cas général, ce n'est pas la convergence uniforme

Comment choisir \mathcal{H} (l'ensemble des hypothèses ?)

1. \mathcal{H} doit être un espace vectoriel
2. $f(x)$ doit avoir un sens...

$$\begin{aligned}\delta_x : \mathcal{H} &\longrightarrow \mathbb{R} \\ f &\longmapsto \delta_x f = f(x)\end{aligned}$$

L^2 ne convient pas !

3. la fonctionnelle d'évaluation doit être continue ($\delta_x \in \mathcal{H}'$) :

$$f_n \xrightarrow{\mathcal{H}} t \implies \forall x \in \mathcal{X}, f_n(x) \xrightarrow{\mathbb{R}} t(x)$$

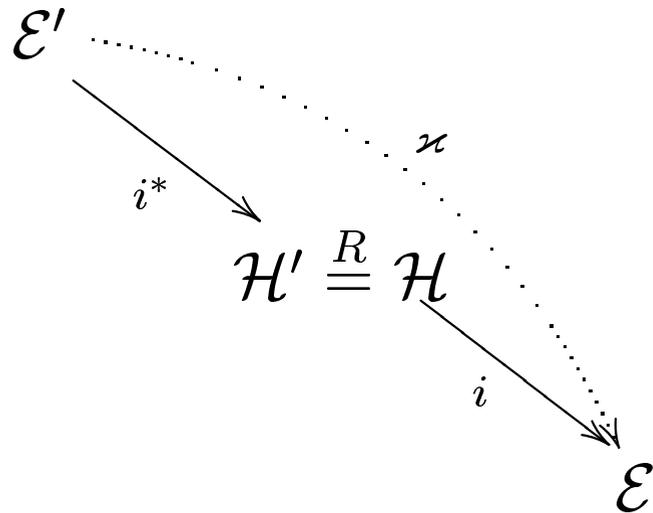
dans le cas général, ce n'est pas la convergence uniforme

(1) + (2) + (3) $\Leftrightarrow \mathcal{H}$ est un espace à noyau reproduisant :
pas nécessairement de Hilbert

\mathcal{H} : un espace à noyau reproduisant

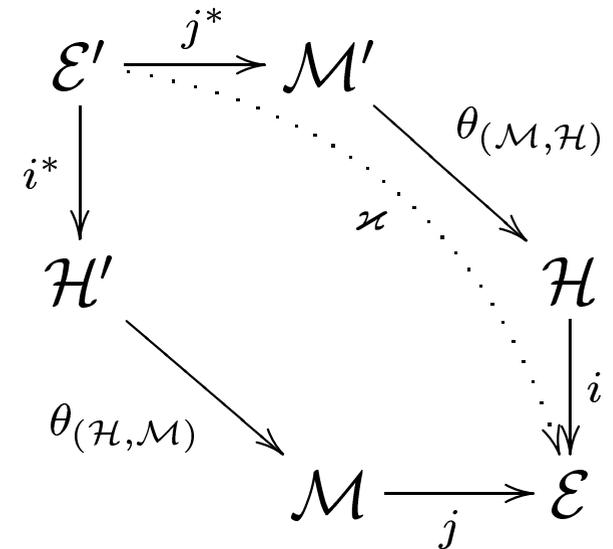
\mathcal{H} espace vectoriel d'évaluation \Rightarrow il existe un noyau $K(s, t)$

Cas Hilbertien



$$K(s, t) = \langle K(s, \cdot), K(\cdot, t) \rangle_{\mathcal{H}}$$

Cas général



$$K(s, t) = (\kappa^*(\delta_s), \kappa(\delta_t))_{\mathcal{M}, \mathcal{H}}$$

mesures

fonctions

Généralisation des noyaux reproduisants et de Schwarz

Construction de \mathcal{H} dans le cas Hilbertien à partir de L^2

- soit S un opérateur linéaire de L^2 sur $\mathbb{R}^{\mathcal{X}}$
- soit $\Gamma(t, \tau)$ une fonction de deux variables, $\Gamma(t, \cdot) = \Gamma_t \in L^2$

$$\begin{aligned} S : L^2 &\longrightarrow \mathbb{R}^{\mathcal{X}} && \text{opérateur de Carleman} \\ f &\longmapsto g(t) = Sf(t) = \int \Gamma(t, \tau) f(\tau) d\tau \end{aligned}$$

Construction de \mathcal{H} dans le cas Hilbertien à partir de L^2

- soit S un opérateur linéaire de L^2 sur $\mathbb{R}^{\mathcal{X}}$
- soit $\Gamma(t, \tau)$ une fonction de deux variables, $\Gamma(t, \cdot) = \Gamma_t \in L^2$

$$\begin{aligned} S : L^2 &\longrightarrow \mathbb{R}^{\mathcal{X}} && \text{opérateur de Carleman} \\ f &\longmapsto g(t) = Sf(t) = \int \Gamma(t, \tau) f(\tau) d\tau \end{aligned}$$

- à partir de S nous définissons $\mathcal{H} = \text{Im}(S)$

Construction de \mathcal{H} dans le cas Hilbertien à partir de L^2

- soit S un opérateur linéaire de L^2 sur $\mathbb{R}^{\mathcal{X}}$
- soit $\Gamma(t, \tau)$ une fonction de deux variables, $\Gamma(t, \cdot) = \Gamma_t \in L^2$

$$\begin{aligned} S : L^2 &\longrightarrow \mathbb{R}^{\mathcal{X}} && \text{opérateur de Carleman} \\ f &\longmapsto g(t) = Sf(t) = \int \Gamma(t, \tau) f(\tau) d\tau \end{aligned}$$

- à partir de S nous définissons $\mathcal{H} = \text{Im}(S)$

$$K(s, t) = \int \Gamma(t, \tau) \Gamma(s, \tau) d\tau = \langle \Gamma(t, \cdot) \Gamma(s, \cdot) \rangle_{L^2} = \langle S\Gamma(t, \cdot) S\Gamma(s, \cdot) \rangle_{\mathcal{H}}$$

$$\|g\|_{\mathcal{H}}^2 = \|Sf\|_{\mathcal{H}}^2 = \|f\|_{L^2}^2 = \|S^{-1}g\|_{L^2}^2$$

Construction de \mathcal{H} dans le cas Hilbertien à partir de L^2

- soit S un opérateur linéaire de L^2 sur $\mathbb{R}^{\mathcal{X}}$
- soit $\Gamma(t, \tau)$ une fonction de deux variables, $\Gamma(t, \cdot) = \Gamma_t \in L^2$

$$\begin{aligned} S : L^2 &\longrightarrow \mathbb{R}^{\mathcal{X}} && \text{opérateur de Carleman} \\ f &\longmapsto g(t) = Sf(t) = \int \Gamma(t, \tau) f(\tau) d\tau \end{aligned}$$

- à partir de S nous définissons $\mathcal{H} = \text{Im}(S)$

$$K(s, t) = \int \Gamma(t, \tau) \Gamma(s, \tau) d\tau = \langle \Gamma(t, \cdot) \Gamma(s, \cdot) \rangle_{L^2} = \langle S\Gamma(t, \cdot) S\Gamma(s, \cdot) \rangle_{\mathcal{H}}$$

$$\|g\|_{\mathcal{H}}^2 = \|Sf\|_{\mathcal{H}}^2 = \|f\|_{L^2}^2 = \|S^{-1}g\|_{L^2}^2 = \|Pf\|_{L^2}^2$$

$$\boxed{S = P^{-1}}$$

Construction de \mathcal{H} dans le cas Hilbertien à partir de L^2

- soit S un opérateur linéaire de L^2 sur $\mathbb{R}^{\mathcal{X}}$
- soit $\Gamma(t, \tau)$ une fonction de deux variables, $\Gamma(t, \cdot) = \Gamma_t \in L^2$

$$\begin{aligned} S : L^2 &\longrightarrow \mathbb{R}^{\mathcal{X}} && \text{opérateur de Carleman} \\ f &\longmapsto g(t) = Sf(t) = \int \Gamma(t, \tau) f(\tau) d\tau \end{aligned}$$

- à partir de S nous définissons $\mathcal{H} = \text{Im}(S)$

$$K(s, t) = \int \Gamma(t, \tau) \Gamma(s, \tau) d\tau = \langle \Gamma(t, \cdot) \Gamma(s, \cdot) \rangle_{L^2} = \langle S\Gamma(t, \cdot) S\Gamma(s, \cdot) \rangle_{\mathcal{H}}$$

$$\|g\|_{\mathcal{H}}^2 = \|Sf\|_{\mathcal{H}}^2 = \|f\|_{L^2}^2 = \|S^{-1}g\|_{L^2}^2 = \|Pf\|_{L^2}^2$$

$$\boxed{S = P^{-1}}$$

L^2 étant séparable \mathcal{H} l'est aussi : il existe une base dénombrable telle que $K(x, y) = \sum \Phi(x)\Phi(y)$

frame - structures obliques

- Famille « interprétable de fonctions » $\{\varphi_k\}_{k \in \mathbb{N}}$
- base + dépendance + redondance
- représentation

$$f = \sum_{k=1}^{\infty} \langle f, \varphi^* \rangle \varphi$$

φ^* : frame dual

- exemple : les ondelettes

bonnes propriétés statistiques - mais une infinité de terme !

construction de frame apprenables (de noyau)

■ Opérateur de Carleman

$$S : L^2 \longrightarrow \mathbb{R}^{\mathcal{X}}$$
$$f \longmapsto Sf = \int \Gamma(\tau, t) f(\tau) d\tau$$

$$- \mathcal{H} = \text{Im}(S)$$

$$K(s, t) = \int \Gamma(s, \tau) \Gamma(\tau, t) d\tau$$

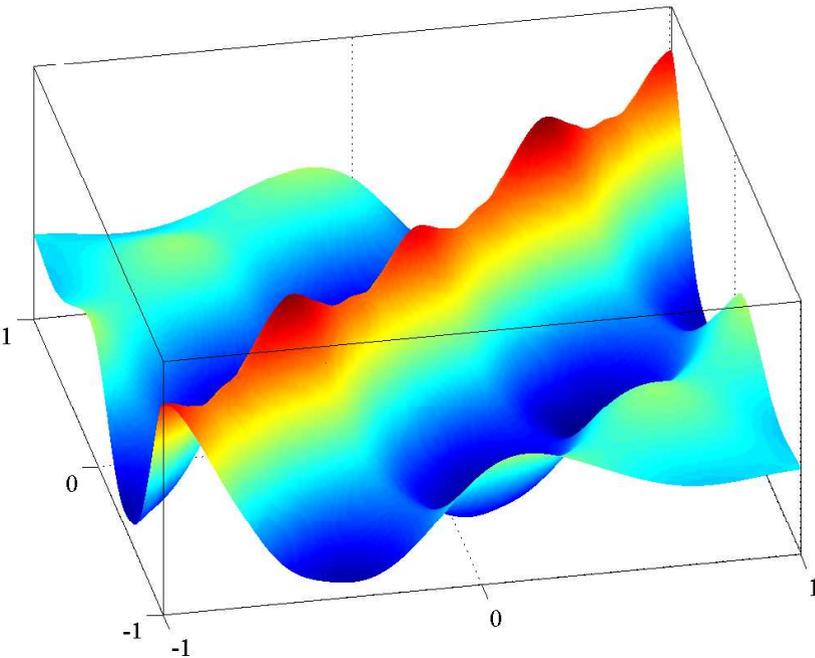
■ Noyaux d'ondelettes

$$- \psi_i \text{ frame de } L^2, \quad \varphi_i = S\psi_i$$

$$- \varphi_i \text{ frame de } \text{Im}(S)$$

$$- K(s, t) = \sum_i \varphi_i^*(s) \varphi_i(t)$$

Meyer Wavelet Kernel $K(x, y)$



La fonction Γ est représentée par une matrice

L'astuce du noyau

Si K est un noyau défini positif d'un k.r.h.s. séparable, il existe une famille $(\phi_j)_{j \in J}$ 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 de Gram.

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

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

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

L'astuce du noyau

Si K est un noyau défini positif d'un k.r.h.s. séparable, il existe une famille $(\phi_j)_{j \in J}$ 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 de Gram.

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

dim ∞

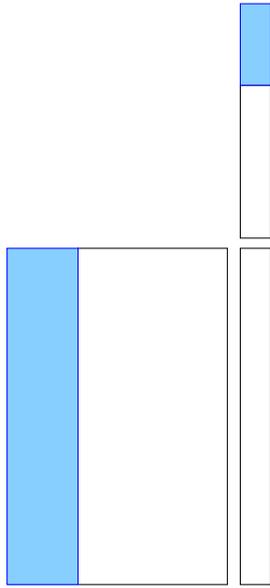
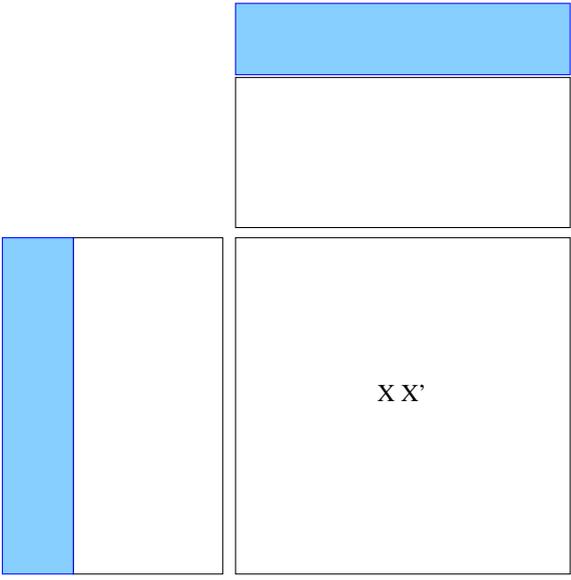
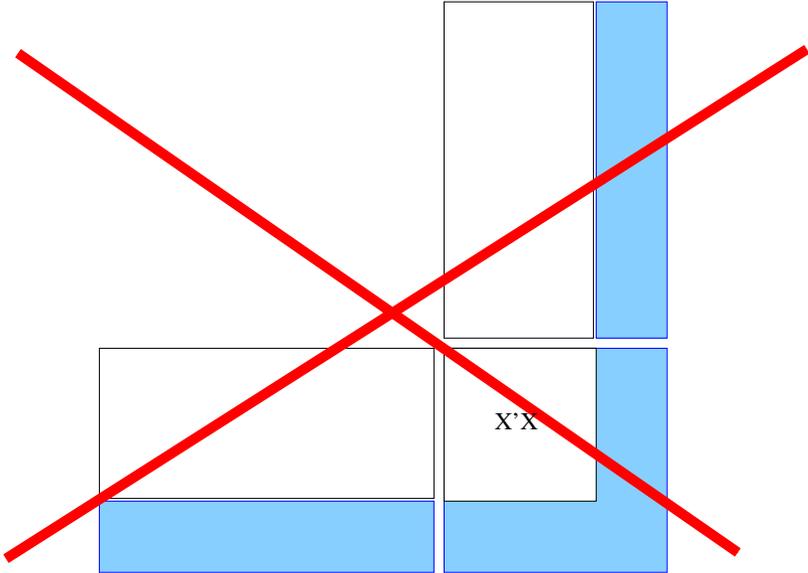
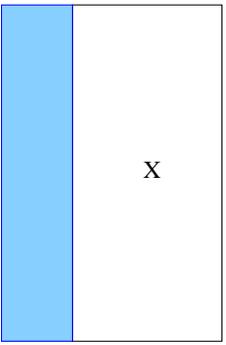
dim n

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

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

Illustration : le cas de l'ACP...

Nouvelles variables...
par exemple x^2 !

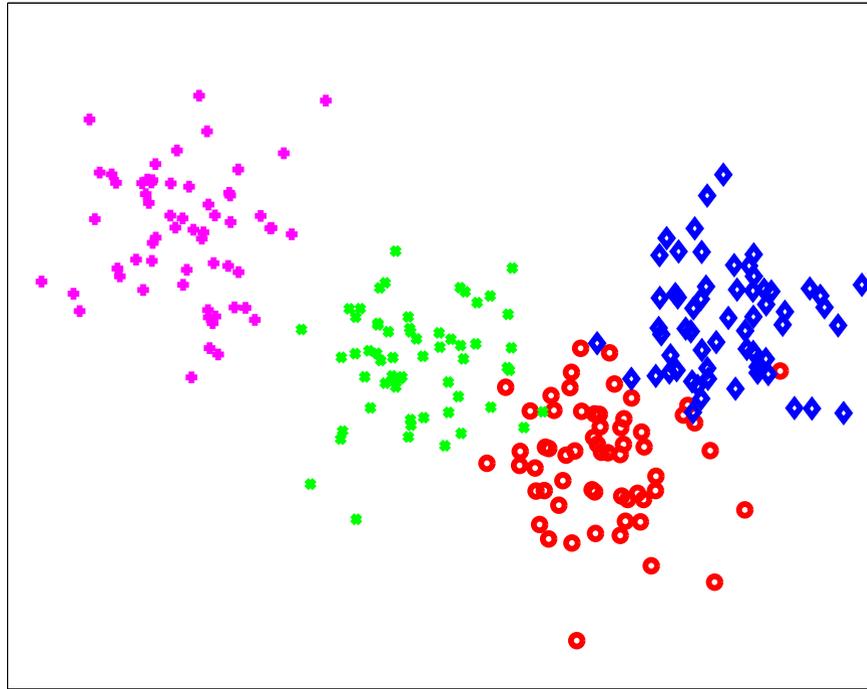


les
variables...

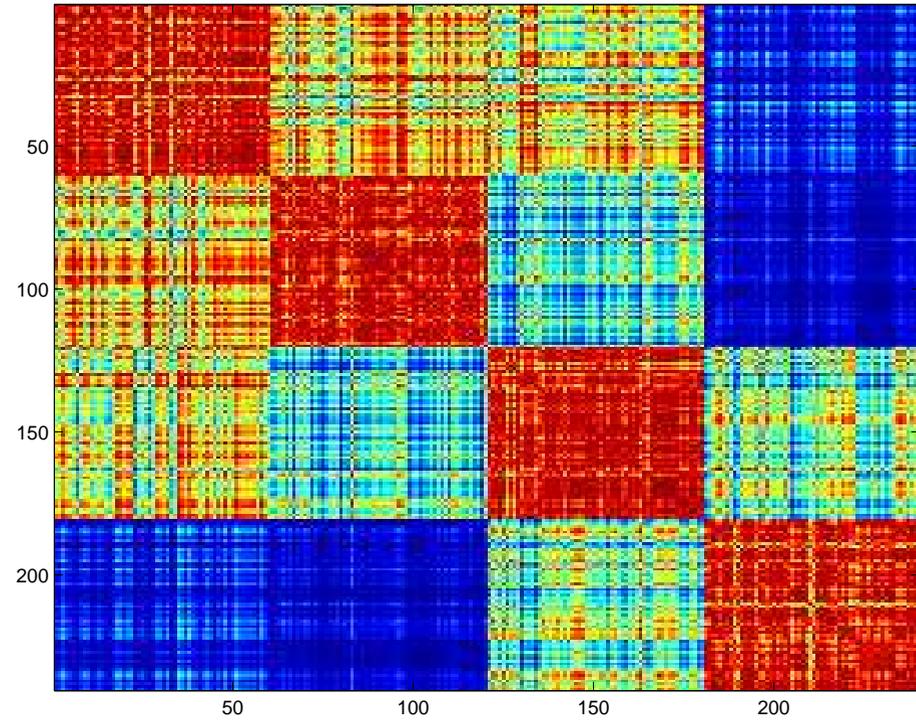
$X^T X$ où XX^T

ou les
individus

La matrice de Gram : illustration



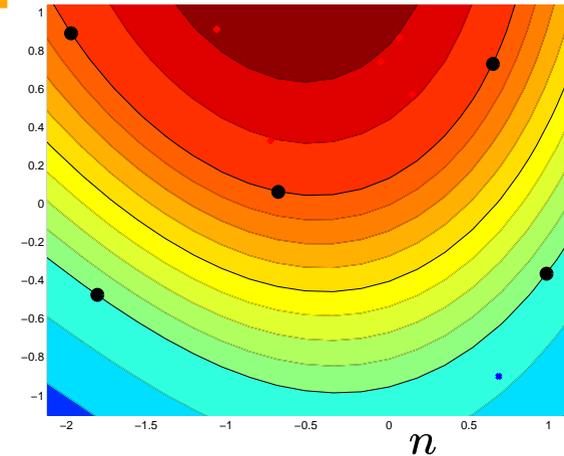
4 classes de 60 points chacune



K matrice de gram associée

$$K_{ij} = \exp - \frac{\|x_i - x_j\|^2}{b}$$

Régularité et marge



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

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

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

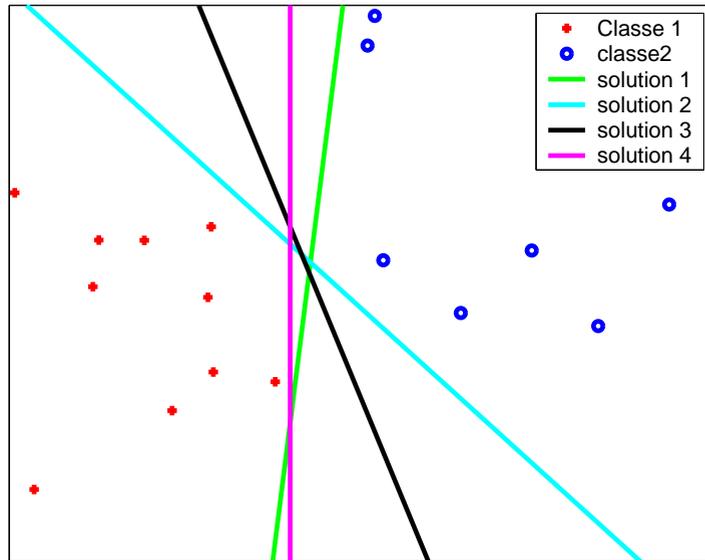
$$\blacksquare \mathbb{P}(err) < \underbrace{\sum_{i=1}^n \mathbb{1}_{\{\text{signe}(f(x_i)) \neq y_i\}}}_{\text{erreur empirique}} + \phi\left(\frac{1}{\text{marge}}\right)$$

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

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

Régularité et marge

Cas linéaire : quelle solution choisir ?



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

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

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

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

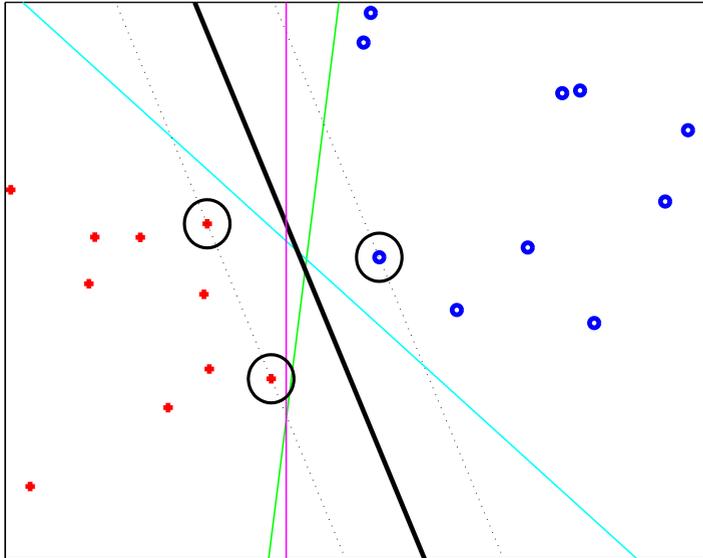
■ minimiser $\mathbb{P}(err) \Leftrightarrow$ maximiser la marge

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

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

Régularité et marge

Celle qui maximise la marge



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

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

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

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

erreur empirique

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

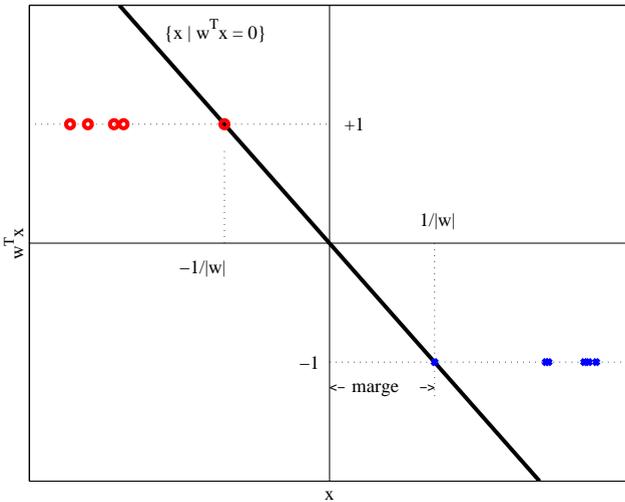
■ maximiser la robustesse \Leftrightarrow maximiser la marge

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

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

Régularité et marge

Valeur de la marge dans le cas monodimensionnel



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

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

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

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

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

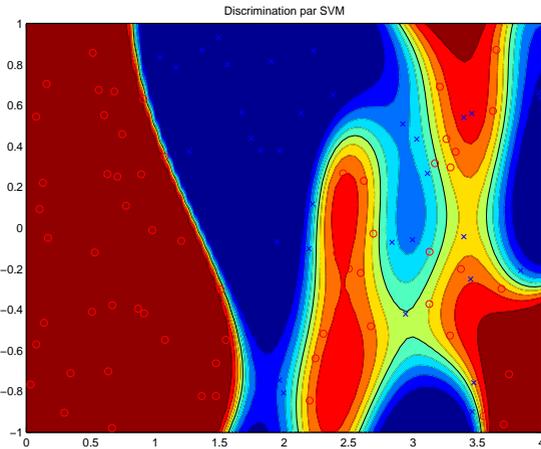
■ maximiser la robustesse \Leftrightarrow maximiser la marge

■ maximiser la marge \Leftrightarrow minimiser $\|\mathbf{w}\|^2$

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

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

Normes...



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

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

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

l'ensemble des hypothèses est alors de la forme

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

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

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 :

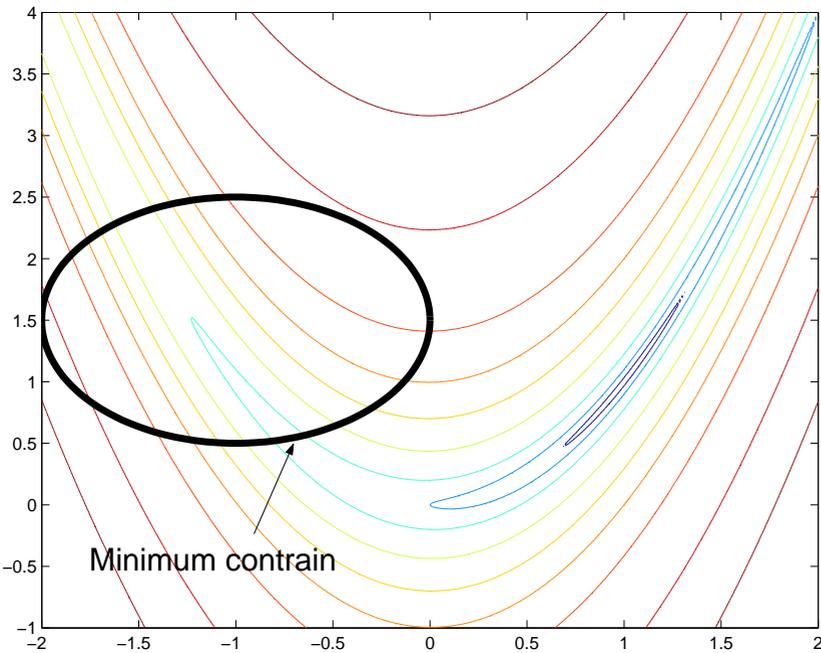
$$\left\{ \begin{array}{l} \min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^{\top} \mathbf{w} + C \sum_{i=1}^n \xi_i \\ \text{avec} \quad y_i f(\mathbf{x}_i) > 1 - \xi_i \quad i = 1, n \\ \text{et} \quad \xi_i > 0 \quad i = 1, n \end{array} \right. \quad (4)$$

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

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

Minimisation sous contraintes (cas séparable)

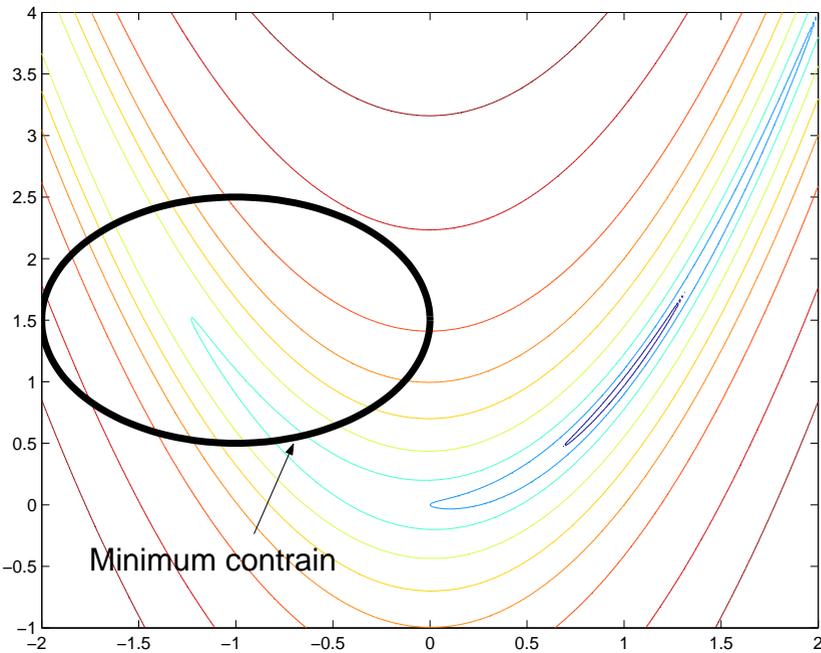
min J(x) dans le domaine admissible



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

Minimisation sous contraintes (cas séparable)

min J(x) dans le domaine admissible



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

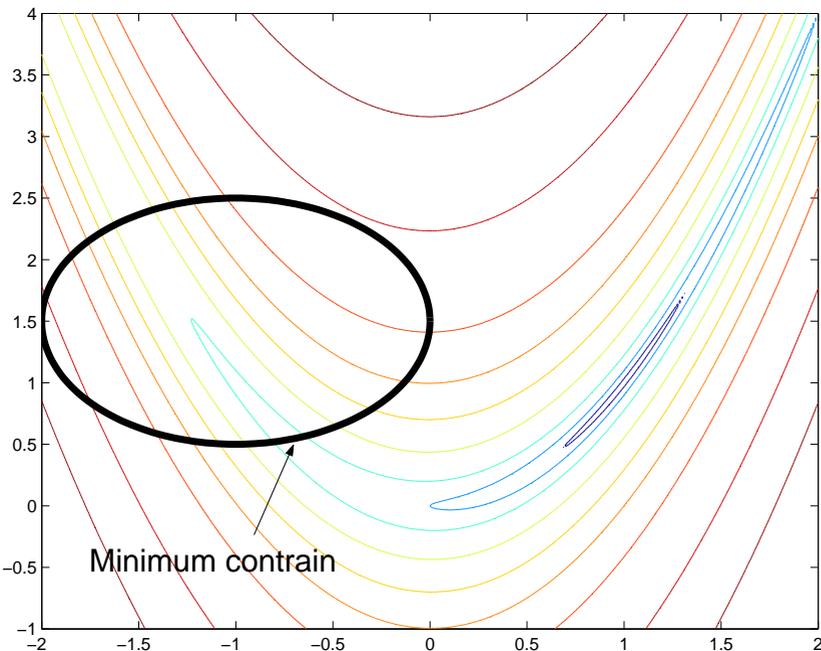
\Leftrightarrow

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

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

Minimisation sous contraintes (cas séparable)

min J(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

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

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

les exemples

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

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

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

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

dont on tire les conditions de Kuhn et Tucker :

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

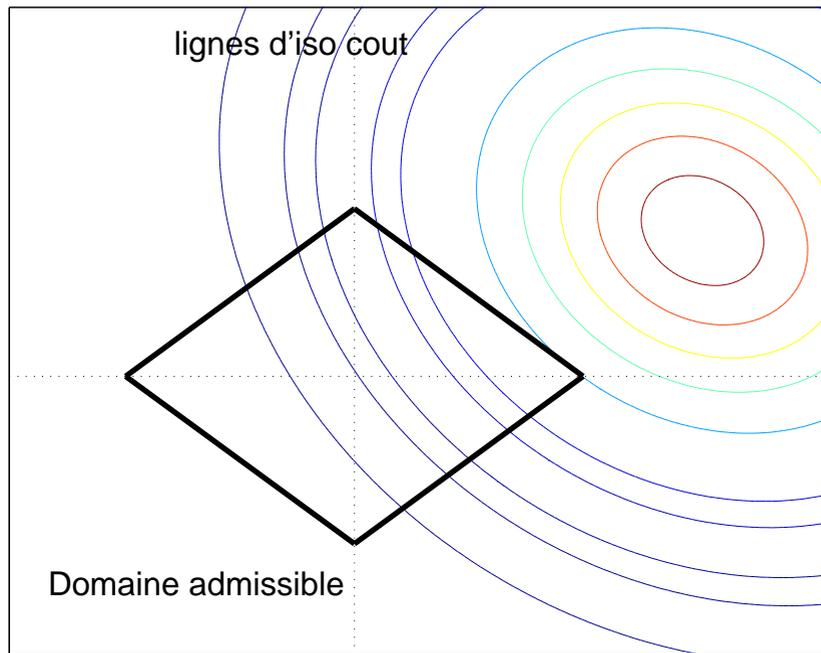
conséquence pour f :

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

Stratégie

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

calcul des λ : problème Dual (2)



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

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.

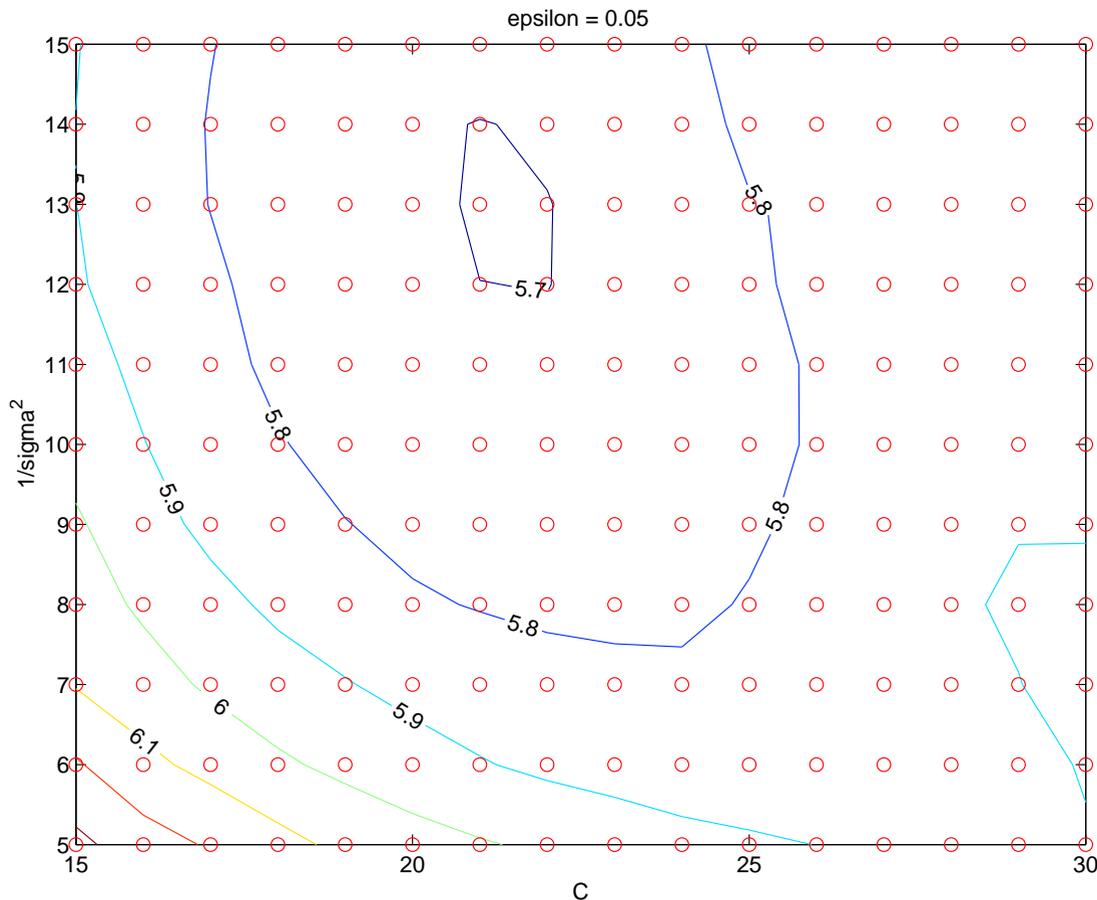
Reformulation de Girosi (97)

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

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

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

Solution pratique : Paramètres de régularisation



■ b : métrique

■ C : des erreurs ?

⇒ la *span bound*

Estimation du taux d'erreur en fonction des hyperparamètres.

SVMtoolbox : une démo ?

```
[xsup, w, w0, tps, alpha] = svmclass(Xapp, yi, C, lambda, kernel, kerneloption, 1, phi);
```

```
ypred = svmval(Xtest, xsup, w, w0, kernel, kerneloption);
```

Conclusion

- SVM les principes
 - Noyaux : universel
 - Marge : minimum global unique
 - Parcimonieux : l'influence de chaque exemple
- SVM et les autres
 - SVM vs Réseaux de neurones (PMC) : optimisation
 - SVM vs Parzen : parcimonie (L^2 vs L^1)
 - SVM : des résultats (vision, génomique, texte)
- les questions sur les SVM
 - le multiclasse
 - comment choisir le noyau ?
comment définir la distance entre deux points ?
 - noyaux non symétrique, non positifs ? d'autres critères ?
 - lissage et détection de ruptures

Références

■ SVM

- Vladimir Vapnik. The Nature of Statistical Learning Theory. Springer, 1995.
- Nello Cristianini and John Shawe-Taylor. An Introduction to Support Vector Machines. Cambridge University Press, Cambridge, UK, 2000
- 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>