

Apprentissage Artificiel

Séparateurs à Vastes Marges

Support Vector Machines

Ewa Kijak

ESIR/Univ. Rennes

Sommaire

Introduction

Les idées de base

Un problème d'optimisation

Redescription en grande dimension

En pratique

Apprentissage de surfaces séparatrices linéaires

- ▶ Dans \mathbb{R}^d , une surface linéaire est un hyperplan H , défini par :

$$w_0 + \mathbf{w}^T \mathbf{x} = 0$$

avec \mathbf{w} vecteur de dimension d et w_0 scalaire.

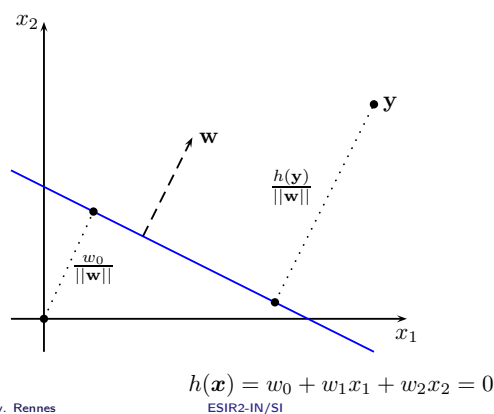
- ▶ Un hyperplan divise l'espace en deux régions, celle où :

$$\mathbf{x} \in \omega_1 \Rightarrow w_0 + \mathbf{w}^T \mathbf{x} \geq 0$$

et l'autre :

$$\mathbf{x} \in \omega_2 \Rightarrow w_0 + \mathbf{w}^T \mathbf{x} \leq 0$$

La géométrie d'un hyperplan



Sommaire

Introduction

Les idées de base

Un problème d'optimisation

Redescription en grande dimension

En pratique

Séparation linéaire dans \mathcal{X}

- ▶ Soit \mathcal{X} , l'espace de représentation des données.
- ▶ Si les points d'apprentissage sont séparables, il existe une infinité d'hyperplans séparateurs.
- ▶ Comment trouver le meilleur hyperplan séparateur dans \mathcal{X} ?
- ▶ Et comment faire si les données ne sont pas séparables ?

Passage dans un autre espace $\Phi(\mathcal{X})$

- ▶ Deux classes peuvent être linéairement séparées dans un espace de grande dimension $\Phi(\mathcal{X})$.
- ▶ Comment trouver le meilleur hyperplan séparateur dans $\Phi(\mathcal{X})$ avec des calculs peu coûteux ?

Séparateurs linéaires et marge

Pour le moment, on se place dans $\mathcal{X} = \mathbb{R}^d$.
 $S = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$, avec $\mathbf{x}^{(i)} \in \mathbb{R}^d$, $y^{(i)} \in \{-1, 1\}$.

On suppose qu'il existe une séparatrice linéaire permettant de distinguer les exemples positifs (supervisés par $y^{(i)} = +1$) des exemples négatifs (supervisés par $y^{(i)} = -1$).

Il existe donc une hypothèse d'apprentissage $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0$ telle que :

$$h(\mathbf{x}^{(i)}) = \mathbf{w}^\top \mathbf{x}^{(i)} + w_0 \begin{cases} > 0 \\ < 0 \end{cases} \implies \hat{y}^{(i)} = \begin{cases} +1 \\ -1 \end{cases}$$

Séparateurs linéaires et marge

Ce qui peut se réécrire :

$$\forall 1 \leq i \leq m : y^{(i)} h(\mathbf{x}^{(i)}) > 0$$

ou

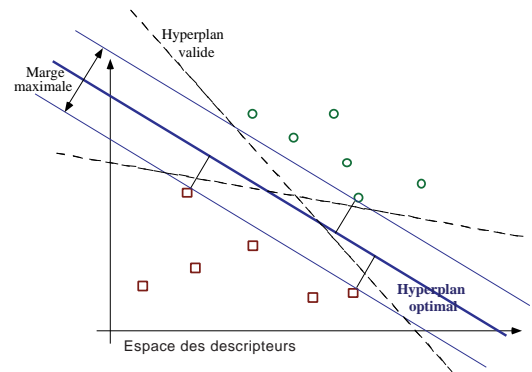
$$\forall 1 \leq i \leq m : y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + w_0) > 0$$

$h(\mathbf{x})$ est l'équation d'un hyperplan dans \mathcal{X} de vecteur normal \mathbf{w} .
 La distance d'un point $\mathbf{x}^{(i)}$ à l'hyperplan d'équation $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + w_0 = 0$ est égale à : $|h(\mathbf{x}^{(i)})| / \|\mathbf{w}\|$.

L'hyperplan est dit sous *forme canonique* lorsque \mathbf{w} et w_0 sont normalisés de façon à satisfaire :

$$\forall 1 \leq i \leq m : y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + w_0) \geq 1$$

Séparateurs linéaires et marge



Le meilleur séparateur linéaire

Lorsqu'il existe une séparatrice linéaire entre les points d'apprentissage, il en existe une infinité. On peut alors chercher parmi ces séparatrices celle qui est la meilleure : la plus écartée des deux nuages de points exemples et contre-exemples.

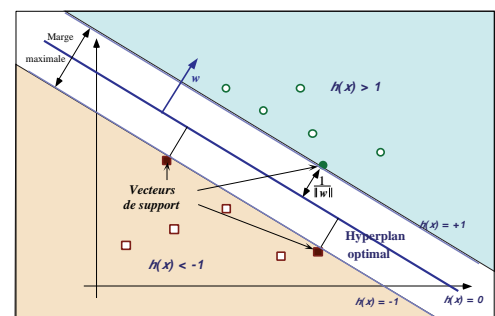
Cet hyperplan optimal est défini par :

$$\operatorname{argmax}_{\mathbf{w}, w_0} \min_i \{ \|\mathbf{x} - \mathbf{x}^{(i)}\| : \mathbf{x} \in \mathbb{R}^d, (\mathbf{w}^\top \mathbf{x} + w_0) = 0, i = 1, \dots, m \}$$

Il maximise la distance minimale aux exemples d'apprentissage.

La *marge* vaut : $2 / \|\mathbf{w}\|$.

Le meilleur séparateur linéaire



L'hyperplan optimal est perpendiculaire au segment de droite le plus court joignant un exemple d'apprentissage à l'hyperplan. Ce segment a pour longueur $\frac{1}{\|\mathbf{w}\|}$ lorsqu'on normalise \mathbf{w} et w_0 .

Sommaire

Introduction

Les idées de base

Un problème d'optimisation

Redescription en grande dimension

En pratique

L'expression primale du problème des SVM

La recherche de l'hyperplan optimal revient donc à minimiser $\|\mathbf{w}\|$, soit à résoudre le problème d'optimisation suivant qui porte sur les paramètres \mathbf{w} et w_0 :

$$\begin{cases} \text{Minimiser} & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{sous les contraintes} & y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + w_0) \geq 1 \quad i = 1, \dots, m \end{cases} \quad (1)$$

Cette écriture du problème, appelée *formulation primale*, implique le réglage de $d + 1$ paramètres, d étant la dimension de l'espace des entrées \mathcal{X} .

- Possible quand d est petit avec méthodes d'optimisation quadratique
- inenvisageable pour des valeurs de d très élevées ($> 10^3$), en particulier quand on travaillera dans $\Phi(\mathcal{X})$.

L'expression duale du problème des SVM

Le *lagrangien* est la somme de la fonction objectif et d'une combinaison linéaire des contraintes. Les coefficients $\alpha_i \geq 0$ sont appelés *multiplicateurs de Lagrange* ou *variables duales*.

$$L(\mathbf{w}, w_0, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y^{(i)} (\mathbf{x}^{(i)} \cdot \mathbf{w} + w_0) - 1) \quad (2)$$

Le problème primal et sa formulation duale ont la même solution qui correspond à un *point-selle* du lagrangien (il faut le minimiser par rapport aux variables primaires \mathbf{w} et w_0 et le maximiser par rapport aux variables duales α_i).

L'expression duale (suite)

Au point-selle, la dérivée du Lagrangien par rapport aux variables primaires doit s'annuler. Ceci s'écrit :

$$\frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, w_0, \alpha) = 0, \quad \frac{\partial}{\partial w_0} L(\mathbf{w}, w_0, \alpha) = 0 \quad (3)$$

et conduit à :

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \quad (4)$$

et à :

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad (5)$$

Solution

En substituant (5) et (4) dans (2), on élimine les variables primaires et l'on obtient la *forme duale* du problème d'optimisation. C'est un problème *quadratique*.

Trouver les multiplicateurs de Lagrange $\alpha_i \geq 0$ tels que :

$$\begin{cases} \max_{\alpha} \left\{ \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}) \right\} \\ \alpha_i \geq 0, \quad i = 1, \dots, m \\ \sum_{i=1}^m \alpha_i y^{(i)} = 0 \end{cases} \quad (6)$$

Intuitif et remarquable

- En pratique, seuls les points qui sont sur les hyperplans frontière $(\mathbf{x}^{(i)} \cdot \mathbf{w}) + w_0 = \pm 1$ jouent un rôle, car les multiplicateurs de Lagrange sont non nuls pour ces seuls points.
- Ils sont appelés **vecteurs de support** ou *exemples critiques*.
- Le vecteur solution \mathbf{w}^* a donc une expression en termes d'un sous-ensemble des exemples d'apprentissage : les exemples critiques $(\mathbf{x}^{(c)}, y^{(c)})$.

$$\mathbf{w}^* = \sum_{i=1}^{m_c} \alpha_c^* y^{(c)} \mathbf{x}^{(c)}$$

- C'est en même temps intuitivement satisfaisant puisque l'on voit bien que l'hyperplan solution est entièrement déterminé par ces exemples.

Solution

L'hyperplan solution correspondant peut alors être écrit :

$$h(\mathbf{x}) = (\mathbf{w}^* \cdot \mathbf{x}) + w_0^* = \sum_{i=1}^m \alpha_i^* y^{(i)} (\mathbf{x} \cdot \mathbf{x}^{(i)}) + w_0^* \quad (7)$$

où les α_i^* sont solution de (6) et w_0 est obtenue en utilisant n'importe quel exemple critique $(\mathbf{x}^{(c)}, y^{(c)})$ dans l'équation :

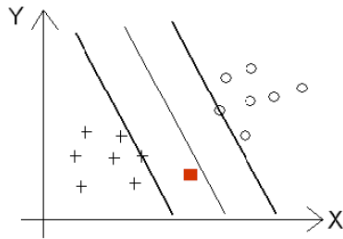
$$y^{(c)} ((\mathbf{x}^{(c)} \cdot \mathbf{w}^*) + w_0) - 1 = 0 \quad (8)$$

Remarques sur la solution

- L'hyperplan solution ne requiert que le calcul des produits scalaires $(\mathbf{x} \cdot \mathbf{x}^{(i)})$ entre des vecteurs de l'espace d'entrée \mathcal{X} .
- La solution ne dépend plus de la dimension d de l'espace d'entrée, mais de la taille m de l'échantillon de données et même du nombre m_c d'exemples critiques qui est généralement bien inférieur à m .
 - nécessite de déterminer m variables duales (α_i) et non $d + 1$ paramètres (\mathbf{w}, w_0) .
- Les méthodes d'optimisation quadratique standard suffisent pour la plupart des cas (env. 10^5 exemples).
- Les fonctions de coût et les contraintes sont strictement convexes (Th. de Kuhn-Tucker)

Classification d'une nouvelle donnée

La classification d'un nouvel exemple inconnu est donnée par sa position par rapport à l'hyperplan optimal.



Dans le schéma ci-dessus, le nouvel élément sera classé dans la catégorie des « + ».

Classification d'une nouvelle donnée

La classe d'une nouvelle donnée \mathbf{x} est fournie par le signe de

$$h(\mathbf{x}) = \mathbf{w}^* \cdot \mathbf{x} + w_0^*$$

$$h(\mathbf{x}) \begin{cases} > 0 \\ < 0 \end{cases} \implies \hat{y} = \begin{cases} +1 \\ -1 \end{cases}$$

On calcule :

$$h(\mathbf{x}) = \sum_{c=1}^{m_c} \alpha_c^* y^{(c)} (\mathbf{x} \cdot \mathbf{x}^{(c)}) + w_0^* \quad (9)$$

où m_c est le nombre de vecteurs support $(\mathbf{x}^{(c)}, y^{(c)})$.

Cas d'un échantillon non linéairement séparable 1

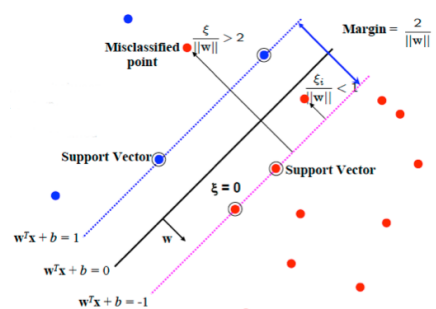
Si les exemples ne peuvent pas être linéairement séparés, on peut employer la technique dite des *variables ressort* (*slack variables*). On modifie les contraintes en les relâchant grâce à des variables ressort $\xi_i \geq 0$:

$$y^{(i)} ((\mathbf{w} \cdot \mathbf{x}^{(i)}) + w_0) \geq 1 - \xi_i \quad (10)$$

Remarque :

- si $\xi_i = 0$, l'exemple $\mathbf{x}^{(i)}$ est bien classé et sur la marge ou à l'extérieur de la marge
- si $0 < \xi_i < 1$, l'exemple $\mathbf{x}^{(i)}$ est bien classé mais à l'intérieur de la marge
- si $\xi_i = 1$, l'exemple $\mathbf{x}^{(i)}$ est sur l'hyperplan optimal (séparatrice)
- si $\xi_i > 1$, l'exemple $\mathbf{x}^{(i)}$ est mal classé

Cas d'un échantillon non linéairement séparable 2



ξ_i est d'autant plus grand que l'exemple $\mathbf{x}^{(i)}$ est loin de la séparatrice.

Echantillon non linéairement séparable 3

Le problème d'optimisation sous contraintes s'écrit alors :

$$\begin{cases} \text{Minimiser} & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \\ \text{sous les contraintes} & y^{(i)} (\mathbf{w}^\top \mathbf{x}^{(i)} + w_0) \geq 1 - \xi_i \text{ et } \xi_i \geq 0 \\ & \forall i = 1, \dots, m \end{cases} \quad (11)$$

- Le coefficient C (appelé constante de régularisation - *trade-off parameter*) règle le compromis entre la marge possible entre les exemples et les erreurs admissibles
 - Remarque : ce n'est pas le nombre de mauvaises classifications qui est minimisé (problème NP-complet), mais la somme des distances aux hyperplans marges

Echantillon non linéairement séparable 4

On obtient alors, comme dans le cas séparable, une formulation duale, avec des contraintes légèrement différentes.

$$\begin{cases} \max_{\alpha} \left\{ \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} (\mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}) \right\} \\ \forall i, \quad 0 \leq \alpha_i \leq C \\ \sum_{i=1}^m \alpha_i y^{(i)} = 0 \end{cases} \quad (12)$$

Les vecteurs support sont toujours tels que $\alpha_i > 0$:

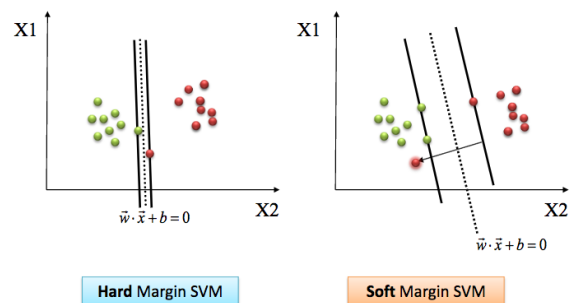
- $0 < \alpha_i < C \Rightarrow \xi_i = 0 \Leftrightarrow \mathbf{x}^{(i)}$ est sur un hyperplan marge
- $\alpha_i = C \Leftrightarrow \xi_i > 0 \Leftrightarrow \mathbf{x}^{(i)}$ est à l'intérieur de la marge (bien ou mal classé)

Echantillon non linéairement séparable 5

- On parle de SVM "à marge molle" (*Soft margin SVM*) ou "marge souple"
 - Plus C est grand, moins on admet d'exemples à l'intérieur de la marge (moins on régularise)
 - C très grande \Leftrightarrow SVM à marge dure
- Les SVM à marges molles ont toujours une solution
- Ils sont plus robustes aux *outliers*
- Les SVM à marges dures n'ont pas d'hyper-paramètre (pas de constante C à déterminer)

Echantillon non linéairement séparable 6

Les vecteurs supports ne sont pas nécessairement sur les marges (\neq SVM à marges dures)



Sommaire

Introduction

Les idées de base

Un problème d'optimisation

Redescription en grande dimension

En pratique

Rappel : Les idées de base

- Certains calculs ne sont pas plus coûteux en transformant l'espace de représentation \mathcal{X} en un espace de grande dimension $\Phi(\mathcal{X})$, par exemple les produits scalaires.
- Deux classes peuvent être linéairement séparées dans un espace de grande dimension.
- Comment trouver le meilleur hyperplan séparateur dans $\Phi(\mathcal{X})$ avec des calculs peu coûteux ?

Le passage par un espace de redescription

Plus la dimension de l'espace de description est grande, plus la probabilité de pouvoir trouver un hyperplan séparateur entre les exemples et les contre-exemples est élevée.

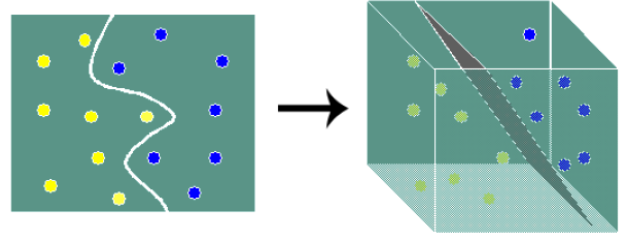
En transformant l'espace d'entrée en un espace de redescription de très grande dimension, éventuellement infinie, il devient donc possible d'envisager d'utiliser la méthode des SVM.

Notons Φ une transformation non linéaire de l'espace d'entrée \mathcal{X} en un espace de redescription $\Phi(\mathcal{X})$:

$$\mathbf{x} = (x_1, \dots, x_d)^\top \xrightarrow{\Phi} \Phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x}), \dots)^\top \quad (13)$$

Le passage par un espace de redescription

$$\begin{aligned} \Phi : \mathbb{R}^d &\mapsto \mathbb{R}^{d'} \quad (d' > d) \\ \mathbf{x} &\mapsto \Phi(\mathbf{x}) \end{aligned}$$



Le passage par un espace de redescription

Le problème d'optimisation se transcrit dans ce cas par :

$$\begin{cases} \max_{\alpha} \left\{ \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \Phi(\mathbf{x}^{(i)}), \Phi(\mathbf{x}^{(j)}) \rangle \right\} \\ \alpha_i \geq 0, \quad i = 1, \dots, m \\ \sum_{i=1}^m \alpha_i y^{(i)} = 0 \end{cases} \quad (14)$$

où $\langle \cdot, \cdot \rangle$ dénote le produit scalaire dans le nouvel espace.

L'équation de l'hyperplan séparateur dans le nouvel espace devient :

$$h(\mathbf{x}) = \sum_{i=1}^m \alpha_i^* y^{(i)} \langle \Phi(\mathbf{x}), \Phi(\mathbf{x}^{(i)}) \rangle + w_0^* \quad (15)$$

où les coefficients α_i^* et w_0^* sont obtenus comme précédemment par résolution de (14).

Le passage par un espace de redescription : les fonctions noyau

$\langle \Phi(\mathbf{x}^{(i)}), \Phi(\mathbf{x}^{(j)}) \rangle$ devient rapidement impossible à calculer quand la dimension de $\Phi(\mathcal{X})$ augmente (sans parler du cas de la dimension infinie), ceci d'autant plus que l'on utilisera des transformations non linéaires des descripteurs d'entrée.

On peut dans certains cas s'arranger pour court-circuiter le passage par les calculs dans l'espace de redescription.

Il existe des fonctions bilinéaires symétriques positives $K(\mathbf{x}, \mathbf{y})$, appelées *fonctions noyau*, faciles à calculer et dont on peut montrer qu'elles correspondent à un produit scalaire $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle$ dans un espace de grande dimension.

Les fonctions noyau : exemple

$$\mathcal{X} \xrightarrow{\Phi} \Phi(\mathcal{X})$$

On connaît une *fonction noyau* $K(\mathbf{x}, \mathbf{y})$, avec \mathbf{x} et \mathbf{y} dans \mathcal{X} telle que

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

Dans ce cas, le produit scalaire dans $\Phi(\mathcal{X})$ est un calcul dans \mathcal{X} .

Par exemple : soient $\mathcal{X} = \mathbb{R}^2$ et $\Phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \in \mathbb{R}^3$

Montrer que : $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle = \langle \mathbf{x}, \mathbf{y} \rangle^2$

Les fonctions noyau

Lorsqu'une telle correspondance est exploitable, le problème d'optimisation (14) est équivalent au problème suivant :

$$\begin{cases} \max_{\alpha} \left\{ \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \right\} \\ \alpha_i \geq 0, \quad i = 1, \dots, m \\ \sum_{i=1}^m \alpha_i y^{(i)} = 0 \end{cases} \quad (16)$$

dont la solution est l'hyperplan séparateur d'équation :

$$h(\mathbf{x}) = \sum_{i=1}^m \alpha_i^* y^{(i)} K(\mathbf{x}, \mathbf{x}^{(i)}) + w_0^* \quad (17)$$

où les coefficients α_i^* et w_0^* sont obtenus comme précédemment par résolution du problème d'optimisation quadratique.

Exemple de fonction noyau : la fonction noyau polynomiale

Il peut être montré que la *fonction noyau polynomiale* :

$$K(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{y} \rangle^n \quad (18)$$

réalise implicitement un produit scalaire dans l'espace des descripteurs correspondant à tous les produits d'exactement n dimensions.

Ainsi pour $n = 2$ et $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$, on a :

$$\langle \mathbf{x}, \mathbf{y} \rangle^2 = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle$$

cf. exemple ci-dessus

qui correspond au changement de description par la fonction :

$$\Phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top \in \mathbb{R}^3$$

La fonction noyau polynomiale

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

$$\Phi : \mathbb{R}^d \mapsto \mathbb{R}^{d'} \quad (d' > d)$$

Pour $n = 2$ et $d = 2$, on peut aussi transformer autrement en dimension $d' = 3$, ou $d' = 4$:

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

$$\Phi(\mathbf{x}) = \left(\frac{1}{\sqrt{2}}(x_1^2 - x_2^2), \sqrt{2}x_1x_2, \frac{1}{\sqrt{2}}(x_1^2 + x_2^2) \right)^\top$$

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

Autre exemple : un polynôme de degré 3

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

$$\Phi : \mathbb{R}^d \mapsto \mathbb{R}^{d'} \quad (d' > d)$$

Pour $n = 3$, $d = 2$ et $d' = 4$, avec :

$$\Phi(\mathbf{x}) = (x_1^3, \sqrt{3}x_1^2x_2, \sqrt{3}x_1x_2^2, x_2^3)^\top$$

on peut également vérifier que :

$$\langle \mathbf{x}, \mathbf{y} \rangle^3 = \langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle$$

Les conditions de Mercer

Une fonction K symétrique est un noyau ssi $(K(x^{(i)}, x^{(j)}))_{i,j}$ est une matrice définie positive.

Dans ce cas, il existe un espace \mathcal{F} et une fonction Φ tels que

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

Si cette condition est vérifiée, on peut appliquer les SVMs

Problèmes :

- Cette condition est très difficile à vérifier
- Elle donne pas d'indication pour la construction de noyaux K
- Elle ne permet pas de savoir comment est Φ

D'autres fonctions noyau

- Noyaux polynomiaux (non homogènes)

$$K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + c)^n$$

- Fonctions à Base Radiale (RBF)

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

- Noyaux gaussiens

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

- Fonctions sigmoïdes

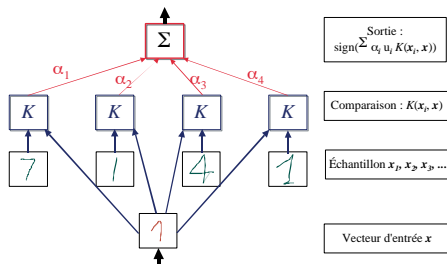
$$K(\mathbf{x}, \mathbf{x}') = \tanh(a(\mathbf{x} \cdot \mathbf{x}') - b)$$

D'autres fonctions noyau

En pratique, on combine des noyaux simples pour en obtenir de plus complexes.

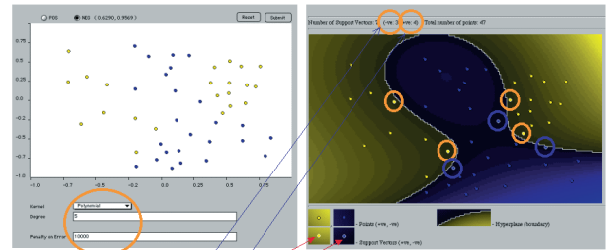
- Construction de nouvelles fonctions noyau par combinaison linéaire de fonctions noyau.

Remarque : Les hyperparamètres (constante de régularisation C , écart-type des gaussiennes σ si on utilise des noyaux gaussiens, degré d si on utilise des noyaux polynomiaux, etc.) doivent être déterminés par validation croisée.



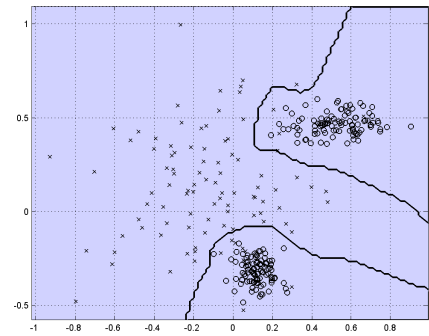
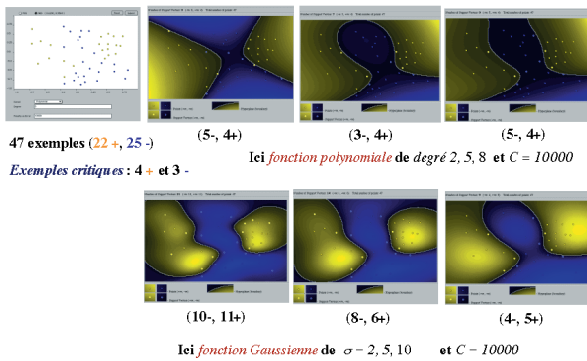
Lors de l'apprentissage, les exemples critiques sont retenus pour définir la fonction de décision. Lorsqu'une nouvelle entrée est présentée au système, elle est comparée aux exemples critiques à l'aide des fonctions noyau qui réalisent un produit scalaire dans l'espace de redescription $\Phi(\mathcal{X})$. La sortie est calculée en faisant une combinaison linéaire de ces comparaisons.

Exemple



- <http://svm.dcs.rhnc.ac.uk/pagones/qpnt.shtml>
- 47 exemples (22 +, 25 -)
- Exemples critiques : 4 + et 3 -
- Ici fonction polynomiale de degré 5 et $C = 10000$

Exemple : Effets de différents choix de la fonction noyau



Séparation de deux classes par SVM avec K fonction à base radiale (RBF).

Sommaire

Introduction

Les idées de base

Un problème d'optimisation

Redescription en grande dimension

En pratique

Implémentation des SVMs

→ Minimisation de fonctions différentiables convexes à plusieurs variables

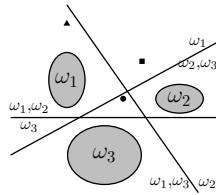
- ▶ pas d'optima locaux
- ▶ mais problèmes de stockage de la matrice noyau (long si milliers d'exemples)
- ⇒ mise au point de méthodes spécifiques
 - ▶ Méthodes itératives, optimisation par morceaux...

Plusieurs librairies publiques disponibles :

- ▶ SVMLight, SVMtorch
- ▶ libSVM, SMO
- ▶ ...

Et pour plus de 2 classes ?

One-versus-all



- ▶ chaque classe est séparée de toutes les autres : il y a C hyperplans
- ▶ exemple : le triangle est assigné à la classe ω_1 , le carré est ambigu entre ω_1 et ω_2 , le point central est ambigu entre les 3 classes

Et pour plus de 2 classes ?

One-versus-all

- ▶ Technique :
 - ▶ pour chaque classe $\{\omega_1, \omega_2, \dots, \omega_C\}$, apprendre un classifieur binaire h_{ω_i}
 - ▶ pour une observation \mathbf{x} :

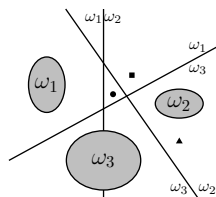
$$\omega^* = \underset{i}{\operatorname{argmax}} h_{\omega_i}(\mathbf{x})$$

→ le classifieur avec la valeur de plus grande confiance (marge) gagne

- ▶ Problème :
 - ▶ calibration : les scores des classifieurs ne sont pas forcément comparables
 - ▶ données non équilibrées : plus d'exemples négatifs que positifs
→ néanmoins : simple et fréquemment utilisé en pratique

Et pour plus de 2 classes ?

One-versus-one



- ▶ les classes sont séparées les unes des autres : il y a $\frac{C(C-1)}{2}$ hyperplans.
- ▶ exemple : le triangle et le carré sont assignés à la classe ω_2 , le point central est ambigu entre les 3 classes

Et pour plus de 2 classes ?

One-versus-one

- ▶ Technique :
 - ▶ pour chaque paire de classes (ω_i, ω_j) , apprendre un classifieur binaire : $\operatorname{sign} h_{ij}$
 - ▶ combiner les classifieurs binaires par un mécanisme de vote majoritaire. Pour une observation \mathbf{x} :

$$\omega^* = \underset{j \in C}{\operatorname{argmax}} |\{i : h_{ij}(\mathbf{x}) = 1\}|$$

→ la classe assignée le plus grand nombre de fois gagne

- ▶ Problème :
 - ▶ computationnel : entraîner $\frac{C(C-1)}{2}$ classifieurs binaires
 - ▶ surapprentissage : la taille de l'ensemble d'apprentissage peut devenir trop petite pour une paire de classes donnée
 - ▶ il peut rester une ambiguïté s'il n'y a pas de classe majoritaire (égalité)

Conclusion

- ▶ Méthode d'apprentissage complètement issue de considérations théoriques (bien fondée mathématiquement).
- ▶ SVM faciles à mettre en oeuvre et donnent souvent de bons résultats en apprentissage (bonne capacité de généralisation : R est proche de R_{emp})
- ▶ ne permettent pas l'extraction d'un modèle compréhensible
- ▶ inadaptés aux très grands volumes de données (calculs lourds, max actuel = 10 000 exemples)
- ▶ pas de solution pour le choix du noyau

Exercice : apprentissage du XOR

- ▶ On souhaite construire un SVM permettant de classer les points $\{(1, 1), (1, -1), (-1, 1), (-1, -1)\}$ selon le résultat de l'opérateur XOR.

Q1 Représenter les points et leur classe dans \mathbb{R}^2 . Ce problème est-il linéairement séparable ?

- ▶ On note Φ une transformation non linéaire de l'espace d'entrée \mathcal{X} en un espace de redescription $\Phi(\mathcal{X})$: $\Phi : \mathbf{x} = (x_1, \dots, x_d) \rightarrow (\phi_1(\mathbf{x}), \dots, \phi_d(\mathbf{x}), \dots)$ et on considère la fonction noyau polynomiale de degré 2 :

$$K(\mathbf{x}, \mathbf{y}) = (\langle \mathbf{x}, \mathbf{y} \rangle + 1)^2$$

Cette fonction réalise implicitement un produit scalaire dans un espace des descripteurs de plus grande dimension : $\langle \Phi(\mathbf{x}), \Phi(\mathbf{y}) \rangle$.

Exercice : apprentissage du XOR

- Q2 Dédurre du développement de $K(\mathbf{x}, \mathbf{y})$ la transformation Φ et la dimension de l'espace de redescription.
- Q3 Rappeler le problème d'optimisation à résoudre faisant intervenir un Lagrangien avec α le vecteur des multiplicateurs de Lagrange.
- Q4 Donner dans un tableau les valeurs $K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ pour tous les couples (i, j) .
- Q5 Ecrire le problème d'optimisation sous la forme d'un système d'équations, et le résoudre.
- Q6 Quel est, dans l'espace $\Phi(\mathcal{X})$, le vecteur poids optimal \mathbf{w}^* ? En déduire l'équation de l'hyperplan optimal.
- Q7 Tracer les séparatrices résultantes dans l'espace d'entrée \mathbb{R}^2 .