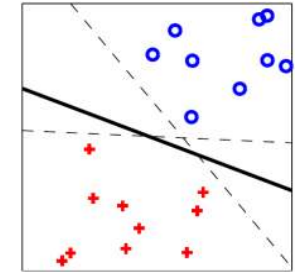


- On étudie tout d'abord un problème bi-classe :
 - Données étiquetées (+ v.s. -)
 - Trouver un hyperplan qui sépare les données de dimension N
 - Hyperplan : espace de dimension $N - 1$
 - En 2d, une droite

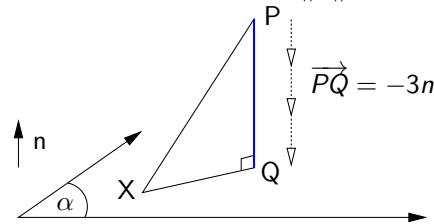


Hypothèse fondamentale

- Données d'apprentissage / de test : distribuées selon la même loi
- En pratique, hypothèse violée à un certain point
- Pour atteindre de bonnes performances, les algorithmes d'apprentissage statistiques doivent disposer d'exemples d'apprentissage suffisamment représentatifs des données de test.

Différentes manières de choisir l'hyperplan car Différents critères à optimiser

- Distance à l'hyperplan séparateur : $d = \frac{|w'x + b|}{\|w\|}$



$d(P, \alpha) = |PQ|$. Le point Q est de la forme $P + tn$, où t est choisi pour que

$$\vec{XQ} \cdot \mathbf{n} = (\vec{XP} + \vec{PQ}) \cdot \mathbf{n} = \vec{XP} \cdot \mathbf{n} + t|\mathbf{n}|^2 = 0.$$

⇒ un seul point et $d(P, \alpha) = |tn| = \frac{|\vec{XP} \cdot \mathbf{n}|}{|\mathbf{n}|}$, (valable aussi lorsque $P \in \alpha$).
Si $\alpha : w_1x_1 + \dots + w_nx_n + b = 0$ (repère orthonormé) :

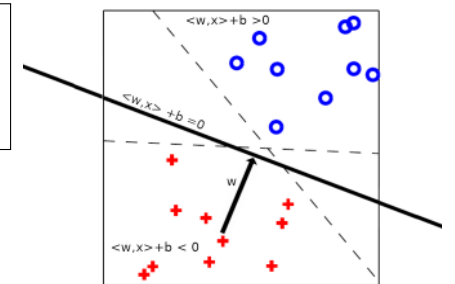
$$d(P, \alpha) = \frac{|w_1p_1 + \dots + w_np_n + b|}{\sqrt{w_1^2 + \dots + w_n^2}}$$

En effet, on peut prendre $\mathbf{n} = (w_1, \dots, w_n)^t$ et si $X = (x_1, \dots, x_n)^t$, on a $\vec{XP} \cdot \mathbf{n} = w_1(p_1 - x_1) + \dots + w_n(p_n - x_n) = w_1p_1 + \dots + w_np_n + b$.

- Distance à l'hyperplan séparateur : $d = \frac{|w'x + b|}{\|w\|}$
 - Critère d'optimisation : minimiser le nombre d'exemples mal classés
- Données mal classées :
- Exemples - : $\langle w, x \rangle + b < 0$
 - Exemples + : $\langle w, x \rangle + b > 0$

Perceptron :

$$\min(- \sum_{\text{mal classés}} y_i (w \cdot x_i + b))$$



Gradient

$$\frac{\partial}{\partial w} = - \sum_{\text{mal classés}} y_i x_i$$

$$\frac{\partial}{\partial b} = - \sum_{\text{mal classés}} y_i$$

Apprentissage

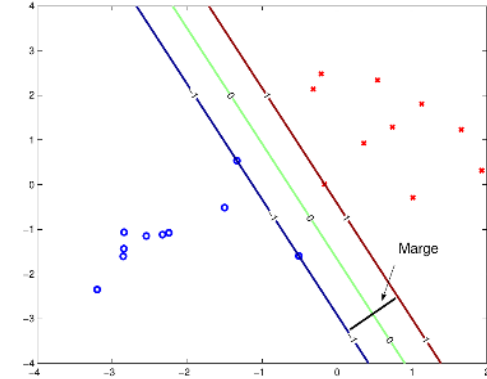
- Initialisation aléatoire des poids
- Mise à jour itérative de l'hyperplan : descente gradient (ρ coeff d'apprentissage)
 - $w \leftarrow w + \rho \sum y_i x_i$
 - $b \leftarrow b + \rho \sum y_i$
- Ou descente gradient stochastique (exemple par exemple)

$$\begin{pmatrix} w \\ b \end{pmatrix}_{n-1} + \rho \begin{pmatrix} y_i \cdot x_i \\ y_i \end{pmatrix} \rightarrow \begin{pmatrix} w \\ b \end{pmatrix}_n$$

Conclusion

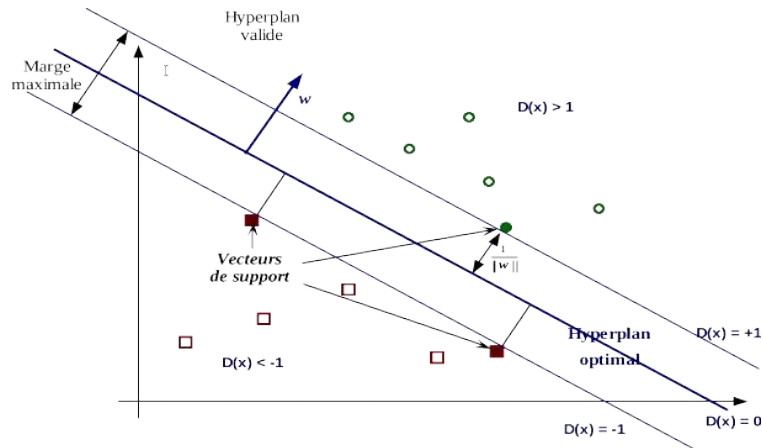
- Supposons les données linéairement séparables : convergence
Pas de convergence sinon
- Solutions multiples au problème dans le cas séparable
Dépendantes de l'initialisation

- Retour sur les séparateurs linéaires
 - Solutions multiples au problème dans le cas séparable
 - Support Vector Machines (SVM) :
Frontière avec « no man's land » maximal, hyperplan « épais »



SVM : maximisation de la marge

- Maximise la "marge" ou rayon du corridor : distance du point le plus proche à l'hyperplan : $1/\|w\|$



SVM : maximisation de la marge

- Maximiser la "marge" : assure de bonnes propriétés de généralisation

$$R < R_{emp} + \sqrt{\frac{1}{n} (h(\ln(2n/h) + 1) - \ln(\alpha/4))} \quad \left(\begin{array}{l} \text{Structural} \\ \text{Risk Minimization} \end{array} \right)$$

- n : nombre d'exemples d'apprentissage
- h : VC dimension ($d + 1$ pour hyperplans dans \mathbb{R}^d)
- Borne valable avec la proba $1 - \alpha$
- Borne sur le risque réel (généralisation)
- Une des raisons principales du succès des SVM
- Maximise la "marge" ou rayon du corridor : distance du point le plus proche à l'hyperplan : $1/\|w\|$
Sous la contrainte que tous les points soient bien classés

$$\begin{cases} \min & \frac{1}{2} \|w\|^2 \\ \forall i & y_i (w \cdot x_i + b) \geq 1 \end{cases}$$

- AI for bio
- J-P Comet
- Introduction
- Trees
- Clustering
- kNN
- Bayes
- SVM
 - Classifieurs linéaires
 - Le perceptron
 - SVM
 - Marge souple
 - Introduction aux noyaux
- NN
- Evaluation
- Data Mining
- Project

- Séparable : \exists hyperplan qui sépare les données sans erreur.
On cherche alors l'hyperplan de marge maximale d'équation :
$$f(x) = \langle w, x \rangle + b = w^t \cdot x + b$$
- Calculons pour X_s un vecteur de support, la distance de X_s à H . Soit P le projeté de X_s sur H et A :
$$\overrightarrow{OP} = \overrightarrow{OX_s} + t \cdot \overrightarrow{w}$$
 où t est choisi telle que $P \in H$:

$$\overrightarrow{w} \cdot \overrightarrow{OP} + b = 0 = \overrightarrow{w} \cdot (\overrightarrow{OX_s} + \overrightarrow{X_s P}) + b = \overrightarrow{w} \cdot (x_s + t \cdot \overrightarrow{w}) + b = 0.$$

Autrement dit, on a : $(\overrightarrow{w} \cdot x_s + b) + t \cdot \|\overrightarrow{w}\|^2 = 0$.

- La marge (plus exactement le double de la marge) est alors donnée par :
$$\text{Marge} = 2d(x, H) = 2 \frac{|w^t \cdot x_s + b|}{\|w\|}$$
- w et b ne sont pas uniques, kw et kb donnent le même hyperplan :

$$kw^t \cdot x + kb = k(w^t \cdot x + b) = 0$$

⇒ Condition de normalisation $|w^t \cdot x_s + b| = 1$ pour les vecteurs de support x_s

ce qui conduit à :
$$\text{Marge} = \frac{2}{\|w\|}$$

- On arrive donc au problème d'optimisation (appelé **problème primal**) :

$$\begin{cases} \min_{w,b} \frac{1}{2} \|w\|^2 \\ \text{tel que } y_i \times (w \cdot x_i + b) \geq 1, i = 1, \dots, n \end{cases}$$

- AI for bio
- J-P Comet
- Introduction
- Trees
- Clustering
- kNN
- Bayes
- SVM
 - Classifieurs linéaires
 - Le perceptron
 - SVM
 - Marge souple
 - Introduction aux noyaux
- NN
- Evaluation
- Data Mining
- Project

- Formulation duale
 - Le dual est un problème quadratique de taille n (# observations).
 - \exists algorithmes bien étudiés et très performants.
 - La formulation duale fait apparaître la matrice de Gram $X \cdot X^t$, ce qui permet de gérer le cas non linéaire à travers des algorithmes à noyaux (qui seront étudiés dans la séance suivante).
- On introduit les multiplicateurs α_i de Lagrange. Le lagrangien est :

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i (w^t \cdot x_i + b) - 1]$$

Le lagrangien doit être optimisé par rapport à w, b et les multiplicateurs α_i .

- En annulant les dérivées partielles du Lagrangien par rapport à w et b , on obtient les relations :

$$\begin{aligned} \frac{\partial}{\partial b} L(w^*, b^*, \alpha^*) = 0 &\Rightarrow \sum_{i=1}^n \alpha_i^* \cdot y_i = 0 \\ \frac{\partial}{\partial w} L(w^*, b^*, \alpha^*) = 0 &\Rightarrow w^* = \sum_{i=1}^n \alpha_i^* \cdot y_i \cdot x_i \end{aligned}$$

- Par substitution de w par $\sum_{i=1}^n \alpha_i^* \cdot y_i \cdot x_i$ dans l'équation du Lagrangien, on obtient le problème dual :

$$\begin{cases} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ t.q \\ \alpha_i \geq 0, i = 1, \dots, n \quad (\text{admissibilité duale}) \\ \sum_{i=1}^n \alpha_i y_i = 0 \quad (\text{stationarité}) \end{cases}$$

- AI for bio
- J-P Comet
- Introduction
- Trees
- Clustering
- kNN
- Bayes
- SVM
 - Classifieurs linéaires
 - Le perceptron
 - SVM
 - Marge souple
 - Introduction aux noyaux
- NN
- Evaluation
- Data Mining
- Project

- Considérons le problème d'optimisation dans \mathbb{R}^n suivant :

$$(P. Primal) \quad \begin{cases} \text{Minimiser} & f(x) \\ & g(x) \leq 0 \\ & x \in S \end{cases}$$

où S est un sous-ensemble de \mathbb{R}^n , $f : S \rightarrow \mathbb{R}$ et $g : S \rightarrow \mathbb{R}^p$.

- Le Lagrangien du problème $L : S \times (\mathbb{R}^+)^p \rightarrow \mathbb{R}$:

$$L(x, u) = f(x) + \langle u, g(x) \rangle \quad \text{où } u \in (\mathbb{R}^+)^p$$

u est appelé le vecteur des multiplicateurs (de Lagrange)

- fonction duale associée au problème $h : (\mathbb{R}^+)^p \rightarrow \mathbb{R}$ telle que :

$$h(u) = \inf_{x \in S} L(x, u) \quad \text{définie sur } U = \{u \in (\mathbb{R}^+)^p \mid \inf_{x \in S} L(x, u) > -\infty\}.$$

- **Théorème** : La fonction h est concave sur tout sous-ensemble convexe de U .
- **Théorème** : $h(u) \leq f(x) \quad \forall u \in U$ et $\forall x \in S$ tel que $g(x) \leq 0$
- Si $\sup_h = +\infty$, le problème primal n'a pas de solution réalisable.
Si $\inf_f = -\infty$, le problème dual n'a pas de solution réalisable.
- Le problème dual associé à (P) est :
(P. Dual) Maximiser $h(u)$ pour $u \in U$

- Introduction
- Trees
- Clustering
- kNN
- Bayes
- SVM
 - Classifieurs linéaires
 - Le perceptron
 - SVM
 - Marge souple
 - Introduction aux noyaux
- NN
- Evaluation
- Data Mining
- Project

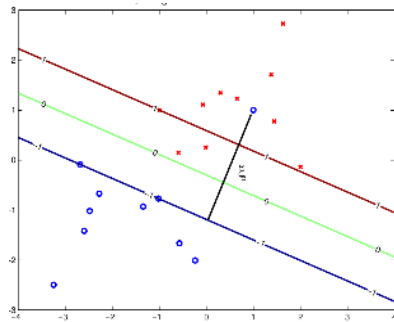
- La solution du problème dual donne les multiplicateurs de Lagrange optimaux α_i^* .
- A partir des α_i
 - on obtient w^* par les relations en haut.
 - Le paramètre b^* est obtenu à partir de la relation $|x_s^T w^* + b^*| = 1$ valable pour tous les vecteurs de support x_s .
- les vecteurs de support sont ceux pour lesquels $\alpha_i \geq 0$. En général leur nombre est beaucoup plus petit que le nombre total d'éléments dans la base d'apprentissage.
⇒ Ajouter des échantillons qui ne sont pas des vecteurs supports à l'ensemble d'apprentissage n'a aucune influence sur la solution finale, c'est à dire seulement les vecteurs de support interviennent dans la fonction de décision (l'expression de la surface séparatrice entre les deux classes).
- Fonction de décision pour une nouvelle observation x :

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

L'hyperplan solution ne dépend que du produit scalaire entre le vecteur d'entrée et les vecteurs de support.

⇒ Cette particularité permet l'utilisation de fonctions noyau pour aborder des problèmes non linéaires (traités dans la séance de cours suivante).

- Les SVM sont des séparateurs linéaires
Que se passe-t-il si on dispose de données d'apprentissage non linéairement séparables ?



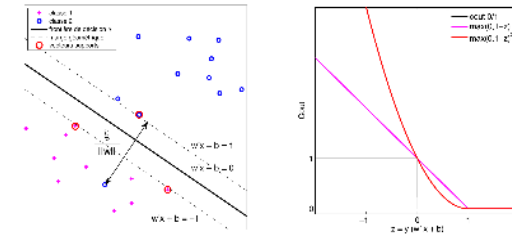
2 solutions :

- Modifier le critère d'apprentissage de manière à autoriser des erreurs d'étiquetage (marge souple)
- Passer dans un espace de représentation où la séparabilité linéaire est possible : noyaux

- L'idée : modéliser les erreurs potentielles par des variables d'écart positives ξ_i associées aux observations $(x_i, y_i), i = 1, \dots, n$.
Si un point (x_i, y_i) vérifie la contrainte de marge $y_i w^T x_i + b \geq 1$ alors la variable d'écart (qui est une mesure du coût de l'erreur) est nulle.

- Nous avons donc deux situations
 - Pas d'erreur : $y_i(w^T x_i + b) \geq 1 \Rightarrow \xi_i = 0$
 - Erreur : $y_i(w^T x_i + b) < 1 \Rightarrow \xi_i = 1 - y_i(w^T x_i + b) > 0$.
 On associe à cette définition une fonction coût appelée "coût charnière" :

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b))$$



Le pb d'optimisation dans le cas des données non-séparables est :

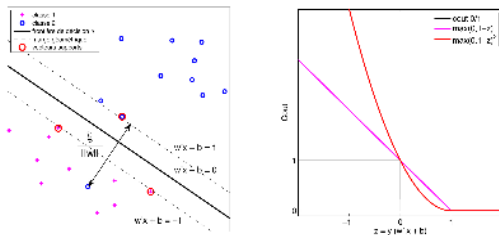
$$\min_{w,b} \left\{ \begin{array}{l} \frac{1}{2} \|w\|^2 \\ \sum_{i=1}^n \xi_i \\ \text{tel que} \begin{cases} y_i(w \cdot x_i + b) \geq 1 - \xi_i, i = 1, \dots, n \\ \xi_i \geq 0, i = 1, \dots, n \end{cases} \end{array} \right.$$

Si $\forall i, \xi_i = 0$, on retrouve le problème linéairement séparable traité plus tôt.

- L'idée : modéliser les erreurs potentielles par des variables d'écart positives ξ_i associées aux observations $(x_i, y_i), i = 1, \dots, n$.
Si un point (x_i, y_i) vérifie la contrainte de marge $y_i w^T x_i + b \geq 1$ alors la variable d'écart (qui est une mesure du coût de l'erreur) est nulle.

- Nous avons donc deux situations
 - Pas d'erreur : $y_i(w^T x_i + b) \geq 1 \Rightarrow \xi_i = 0$
 - Erreur : $y_i(w^T x_i + b) < 1 \Rightarrow \xi_i = 1 - y_i(w^T x_i + b) > 0$.
 On associe à cette définition une fonction coût appelée "coût charnière" :

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b))$$



Le pb d'optimisation dans le cas des données non-séparables est :

$$\min_{w,b} \left\{ \begin{array}{l} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i \\ \text{tel que} \begin{cases} y_i(w \cdot x_i + b) \geq 1 - \xi_i, i = 1, \dots, n \\ \xi_i \geq 0, i = 1, \dots, n \end{cases} \end{array} \right.$$

Si $\forall i, \xi_i = 0$, on retrouve le problème linéairement séparable traité plus tôt.

- On introduit les multiplicateurs α_i de Lagrange. Le lagrangien est :

$$L_p(w, b, \alpha) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [y_i(w^T \cdot x_i + b) - 1 + \xi_i] - \sum_{i=1}^n \mu_i \xi_i$$

Le lagrangien doit être optimisé par rapport à w, b et les multiplicateurs α_i et μ_i .

- En annulant les dérivées partielles du Lagrangien par rapport à w, b et ξ_i , on obtient les relations :

$$\begin{aligned} \frac{\partial L}{\partial b} L(w^*, b^*, \alpha^*) = 0 &\Rightarrow \sum_{i=1}^n \alpha_i^* \cdot y_i = 0 \\ \frac{\partial L}{\partial w} L(w^*, b^*, \alpha^*) = 0 &\Rightarrow w^* = \sum_{i=1}^n \alpha_i^* \cdot y_i \cdot x_i \\ \frac{\partial L}{\partial \xi_i} L(w^*, b^*, \alpha^*) = 0 &\Rightarrow \mu_i = C - \alpha_i \end{aligned}$$

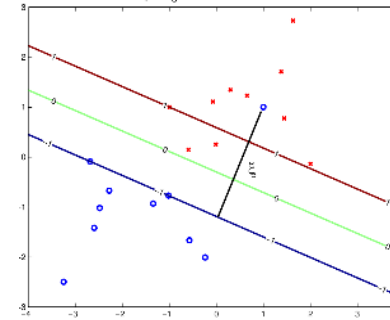
- Par substitution de w par $\sum_{i=1}^n \alpha_i^* \cdot y_i \cdot x_i$ dans l'équation du Lagrangien, on obtient le problème dual :

$$(P. Dual) \left\{ \begin{array}{l} \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ t. q \\ C \geq \alpha_i \geq 0, i = 1, \dots, n \text{ (admissibilité duale)} \\ \sum_{i=1}^n \alpha_i y_i = 0 \text{ (stationarité)} \end{array} \right.$$

- Observations :
 - C joue le rôle d'une constante de régularisation (la régularisation est d'autant plus forte que C est proche de 0).
la régularisation : processus consistant à ajouter de l'information à un problème pour éviter le surapprentissage.
 - généralement, pénalité envers la complexité du modèle.
 - On peut relier cette méthode au principe du rasoir d'Occam.
 - D'un point de vue bayésien, l'utilisation de la régularisation revient à imposer une distribution a priori sur les paramètres du modèle.
 - La différence pour le problème dual entre le cas séparable et non séparable est que les valeurs des α_i sont majorées par C.
 - Les points mal classés ou placés dans la marge ont un $\alpha_i = C$.
 - b est calculé de sorte que $y_i f(x_i) = 1$ pour les points tels que $C > \alpha_i > 0$.
- La fonction de décision permettant de classer une nouvelle observation x est toujours

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

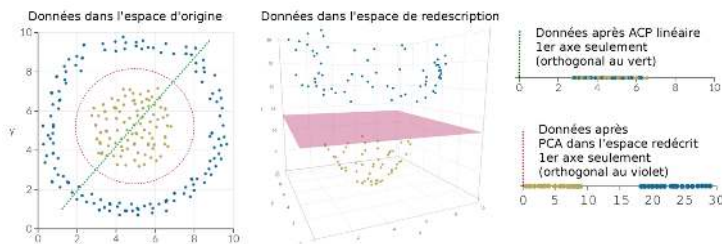
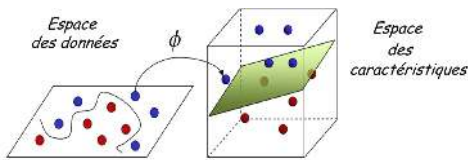
- Les SVM sont des séparateurs linéaires
Que se passe-t-il si on dispose de données d'apprentissage non linéairement séparables ?



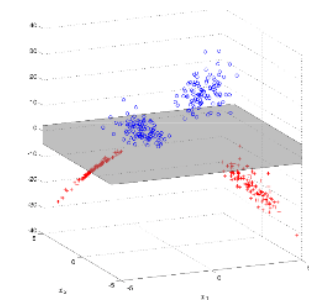
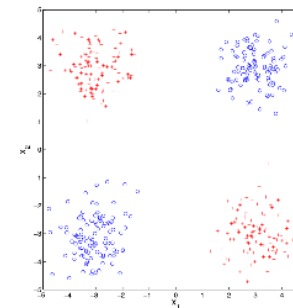
2 solutions :

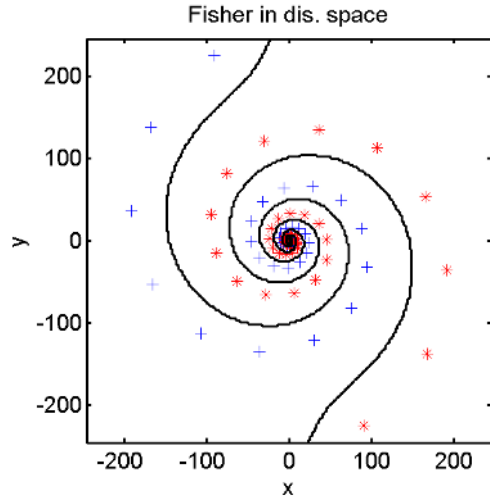
- Modifier le critère d'apprentissage de manière à autoriser des erreurs d'étiquetage (marge souple)
- Passer dans un espace de représentation où la séparabilité linéaire est possible : noyaux

- Autre possibilité pour surmonter le problème des données non linéairement séparables dans l'espace d'entrée (input space) : Passer dans un nouvelle espace Φ (feature space) de grande dimension
- Un séparateur linéaire dans $\Phi(E)$ donne un séparateur non-linéaire dans E.



- input space : (x_1, x_2) , en 2D
- feature space : $(x_1, x_2, x_1 x_2)$, en 3D





c'est un hyperplan, dans un certain espace...

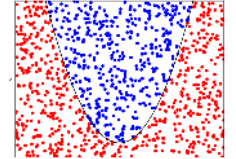
- Solution du SVM dans l'espace induit (feature space)

$$f(x) = \sum_{i=1}^n \alpha_i y_i \langle \Phi(x_i), \Phi(x) \rangle + b$$

- Au lieu de définir explicitement Φ , on préfère définir K

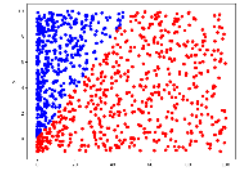
$$K(x, x') = \langle \Phi(x); \Phi(x') \rangle$$

⇒ permet de faire des calculs dans l'espace de départ
Utile, surtout si $dim(\Phi) = \infty$



Exemple :

$$\begin{aligned} x &= (x_1, x_2) \\ \Phi(x) &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \\ \Phi(x)\Phi(x') &= x_1^2x_1'^2 + 2x_1x_2x_1'x_2' + x_2^2x_2'^2 \\ &= (x_1x_1' + x_2x_2')^2 = (xx')^2 \\ \Rightarrow K(x, x') &= (xx')^2 \end{aligned}$$



- On peut donc exprimer la solution SVM sans expliciter Φ

$$f(x) = \sum_{i \in \text{support}} \alpha_i y_i K(x_i, x) + b$$

- On peut choisir n'importe quelle fonction K
Pourvu qu'on puisse prouver qu'il existe un espace dans lequel $K(x, x') = \langle \Phi(x), \Phi(x') \rangle$
- Condition de Mercer :
 $k(x_i, x_j)$ terme général d'une matrice Semi Définie Positive (SDP)
Valeur propres ≥ 0

- Linéaire : $K(x, x') = \langle x|x' \rangle$
- Polynomiale : $K(x, x') = (\langle x|x' \rangle)^d$ ou $(\langle x|x' \rangle + 1)^d$
- Gaussienne (radial basis) : $K(x, x') = e^{-\frac{1}{2\sigma^2} \|x - x'\|^2}$

- Intuitivement, une fonction noyau est proche d'une mesure de similarité. Si x et y sont 2 vecteurs, alors un noyau défini pour ces vecteurs doit avoir des valeurs élevées si $x \approx y$ et des valeurs faibles si x et y sont très différents.
- Une fonction $s : X \times X \rightarrow \mathbb{R}$ est une mesure de similarité si :
 - $\forall x, y \in X, s(x, y) \geq 0$ (positivité)
 - $\forall x, y \in X, s(x, y) = s(y, x)$ (symétrie)
 - $\forall y \in X, y \neq x, s(x, y) < s(x, x)$ (uniformité)
 - $s(x, y) = s(x, x) \Leftrightarrow x = y$ (identité)

A comparer avec le produit scalaire : deux vecteurs proches (directions presque parallèles) ont une valeur élevée pour le produit scalaire.

- Théorème de Mercer. Soit X un compact dans \mathbb{R}^d (compact = fermé et borné) et $K : X \times X \rightarrow \mathbb{R}$ une fonction symétrique. On suppose aussi que $\forall f \in L^2(X) :$

$$\sum_X K(x, y) f(x) f(y) dx dy \geq 0 \quad (\text{condition de Mercer})$$

Alors il existe un espace de Hilbert H et $\Phi : X \rightarrow H$ tel que $\forall x, y \in X :$

$$K(x, y) = \langle \Phi(x), \Phi(y) \rangle \quad (\text{produit scalaire})$$

- La fonction $K(x, y)$ s'appelle noyau positif défini.
- Une condition équivalente pour que la fonction $K : X \times X \rightarrow \mathbb{R}$ soit un noyau positif défini est la suivante :

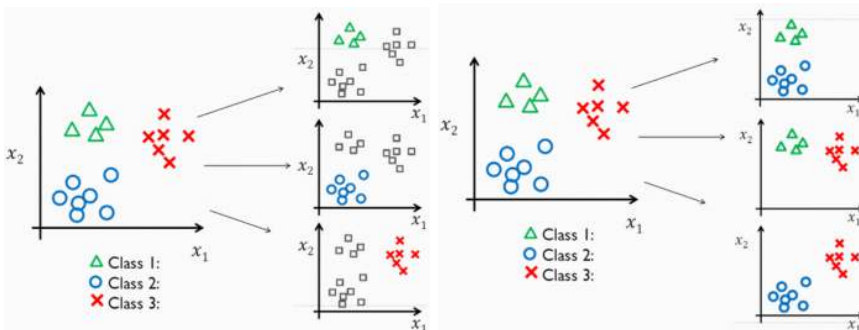
$$\forall n \in \mathbb{N} \text{ et } \{x_i\}_{i=1, \dots, n} \subset X$$

la matrice de Gramm $K = [K_{i,j}]_{i,j=1, \dots, n} = [K(x_i, x_j)]_{i,j=1, \dots, n}$

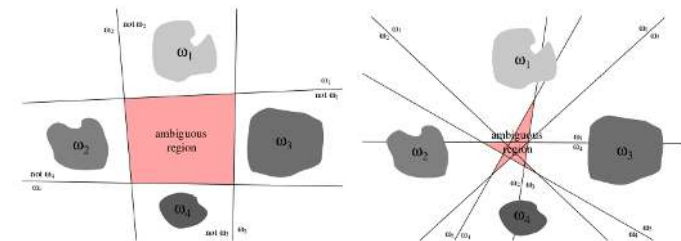
est définie positive, c'est à dire : $\forall c \in \mathbb{R}^n, c \neq 0, \text{ on a } c^T K c > 0$

- Un noyau valide garantit donc l'existence de H et peut s'exprimer alors comme un produit scalaire dans H . Un noyau valide garantit aussi la convexité du problème d'optimisation quadratique sous contraintes d'inégalité rencontré pour les SVM.

- Solutions pour les problèmes multi-catégoriels :
Convertir le problème en un ensemble de problèmes bi-classes
- "one-versus-rest" / "one-against-all"
Pour chaque classe C_i , on recherche une fonction discriminante qui sépare les exemples de la classes C_i de tous les autres exemples
- "one-versus-one"
 $c(c-1)/2$ fonctions discriminantes sont utilisés, une par paire de classes



- Convertir un problème multi-classes en un ensemble de problèmes biclasses peut mener à des régions dans lesquelles la classification est indéfinies...



- OvR : on choisit la classe prédite avec le meilleur score
- OvO : on choisit la classe qui remporte le plus de duels.
- ♥ Avantage d'OvO : taille des ensembles d'apprentissage faible, très utile pour SVM car SVM ne passe pas bien à l'échelle