

## Chap4: SVM et introduction aux classification à noyaux

### I) Introduction

- **SVM: Support vector machine**

Séparation à vaste Marge

- Classifieur binaire, paramétrique, originalement linéaire mais permettant d'en dériver facilement des classifieurs non-linéaire (astuces du noyau)
- Intrinsèquement lié à la fonction de coût de substitution de hinge
- Performante et très utilisé en pratique

#### CHEMINEMENT DU COURS:

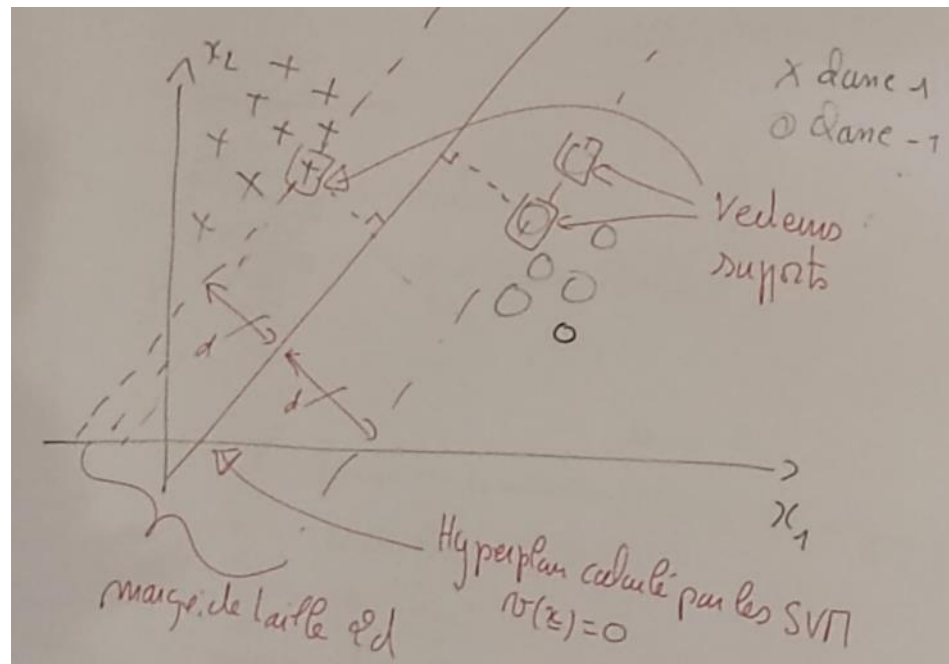
- Cas linéairement séparable
  - Notion de marge (dure)
  - Ecriture comme un pb d'optimisation quadratique sous contrainte.
    - Formulation primale
    - Formulation duale
  - Lien avec hinge
- Cas non linéairement séparable
  - Notion de marge souple
  - Modification de la forme duale
- Classifieur non LINEAIRE
  - Astuce du noyau
  - Lié à forme dual du problème linéaire

### II) SVM LINEAIRE: CAS LINEAIREMENT SEPARABLE

#### 1. Principe

On suppose donc ici que l'on peut trouver des classifieurs linéaires ne faisant aucune erreur sur les données d'apprentissage.

Le principe des SVM est de chercher l'hyperplan séparateur qui offre la plus grande marge de sécurité pour la phase de prédiction. C'est à dire correspondant à la plus grande distance entre l'hyperplan séparateur et les individus qui en sont le plus proche.



Pour simplifier  $y = \{-1, 1\}$ , et on a tjrs

$$v(x) = w_0 + \sum_{n=1}^N w_n x_n \quad w_0 \text{ est le biais}$$

$$v(x) = w_0 + W^T x = \theta^T \phi_x \quad \text{avec } \theta = [w_0 \dots w_n]^T \text{ et } \phi_{x_{vect}} = [1, x_1 \dots x_N]^T$$

$$\text{Et donc} \quad f(x) = \text{sign}(w_0 + W^T x) = \text{sign}(\theta^T \phi_x)$$

## 2. Forme Primale

Notons  $(x_m, y_m)_{\{m=1 \dots M\}}$  la base d'apprentissage servant à déterminer  $w_0$  et  $w^*$

Rappel: La distance  $d_m$  entre  $x_m$  et l'hyperplan  $v(x) = 0$  et donné par :  $d_m = \frac{v(x_m)}{\|w\|}$

Cette distance est "orientée" ( $d_m$  peut être négatif)

On va donc chercher  $(w_0, w^*)$  qui maximise la plus petite distance observable  $|d_m|$  sur l'ensemble des individus:

En remarquant que :

$$|d_m| = \frac{y_m v(x_m)}{\|w\|}$$

Si le classifieur ne fait pas d'erreur.

$$\Rightarrow (w_0, w^*) = \text{argmax} \left\{ \frac{1}{\|w\|} \min \{ y_m (w_0 + w^T x_m) \} \right\} \quad (1)$$

Tel quel le pb semble compliqué mais on peut remarquer que si  $v(x) = 0$  alors  $\beta v(x) = 0 \quad \beta \neq 0$

Donc si  $(w_0, w^*)$  est solution  $(\beta w_0, \beta w^*)$  l'est aussi.

On peut donc chercher comme solution à déterminer le couple  $(w_0, w^*)$  tq  $\min \{ y_m (w_0 + w^T x_m) \} = 1$

Cela revient à imposer que la distance entre l'hyperplan  $(w_0, w^*)$  et l'individu le plus proche vaut  $\frac{1}{\|w\|}$

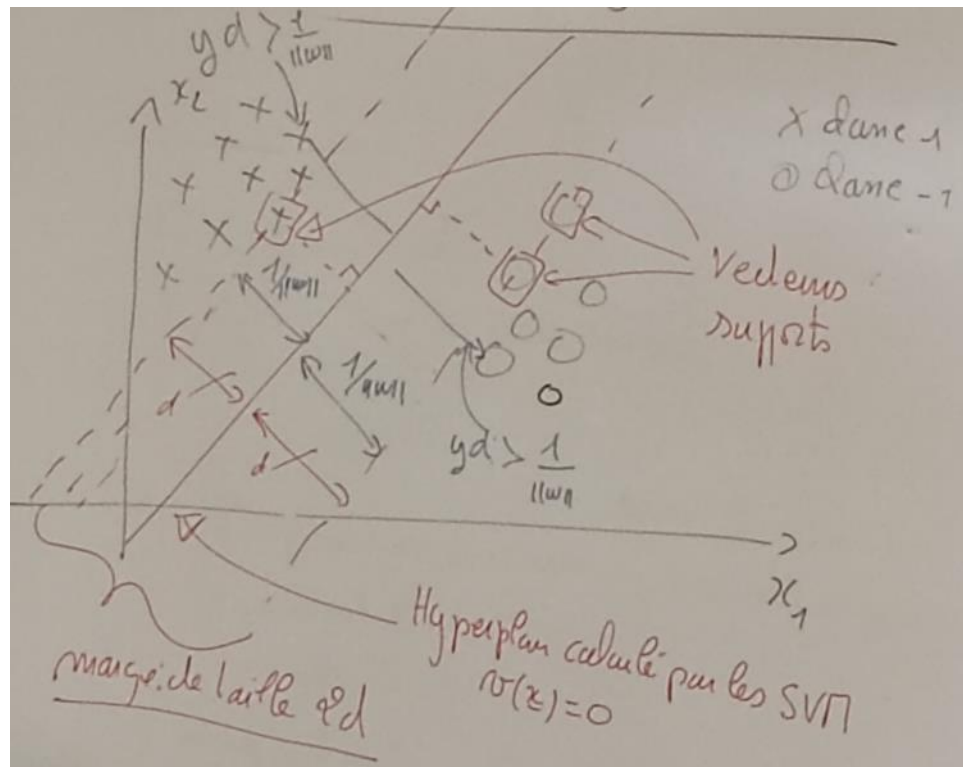
Autrement dit la marge est de taille  $\frac{2}{\|w\|}$

Les Vecteurs supports vérifiant  $yv(x) = 1$ .

Le pb revient donc:

$$(\mathbf{w}_0, \mathbf{w}^*) = \underset{\mathbf{w}}{\operatorname{argmax}} \left\{ \frac{1}{\|\mathbf{w}\|} \right\} \rightarrow \text{maximise la marge}$$

Sous les contraintes: qqs  $m, \mathbf{y}_m(\mathbf{w}_0 + \mathbf{w}^T \mathbf{x}_m) \geq 1 \rightarrow (\text{l'hyperplan est séparateur})$



En pratique, on résout le pb équivalent suivant:

$$(\mathbf{w}_0, \mathbf{w}^*) = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{1}{2} * \|\mathbf{W}\|^2 \right\}$$

Sous les contraintes: qqs  $m, \mathbf{y}_m(\mathbf{w}_0 + \mathbf{w}^T \mathbf{x}_m) \geq 1$  :

**FORME PRIMALE DU PB DES SVM**

Rem:  $\|\mathbf{W}\|^2 = \mathbf{W}^T \mathbf{I}_N \mathbf{W} = \text{forme quadratique def positive}$

Pb classique de minimisation d'une forme quadratique def pos sous contraintes linéaires (type  $\mathbf{A}\boldsymbol{\theta} - \mathbf{b} \geq \mathbf{0}$ )

$\boldsymbol{\theta}$  variable à optimiser (programmation quadratique)

Il existe de nombreux algorithmes efficaces pour résoudre ce pb.

En optimisation sous contraintes, on cherche souvent à réécrire le pb original (primal) sous une autre forme (**duale**) en faisant, impliquant de nouvelles variables.

On introduit ainsi  $M$  variables auxiliaires :

$$\boldsymbol{\alpha}_1 \dots \boldsymbol{\alpha}_M \quad (\boldsymbol{\alpha} = [\boldsymbol{\alpha}_1 \dots \boldsymbol{\alpha}_M]^T)$$

Afin de former le Lagrangien du pb primal

$$\text{Lag}(\mathbf{w}, \mathbf{w}_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{W}\|^2 - \sum_{m=1}^M \alpha_m (\mathbf{y}_m(\mathbf{w}_0 + \mathbf{w}^T \mathbf{x}_m) - 1)$$

On peut montrer alors que le pb primal est équivalent à déterminer  $(\mathbf{w}^*, \mathbf{w}_0^*, \boldsymbol{\alpha}^*)$  tq:

$$\bullet \quad \frac{\partial \text{Lag}(\mathbf{w}, \mathbf{w}_0, \boldsymbol{\alpha})}{\partial \mathbf{w}} = 0 \quad (\alpha)$$

- $\frac{\partial \text{Lag}(w_0)}{\partial w} = 0$  (b)
- Qqs  $\alpha$ ,  $\text{Lag}(w, w_0, \alpha) \leq \text{Lag}(w, w_0, \alpha^*)$  *Lag est max en  $\alpha^*$*   
Sous contraintes qqs m,  $\{\alpha_m(y_m(w_0 + w^T x_m) - 1) = 0$  (e1) et  $\alpha_m \geq 0$  (e2)}

(a) et (b) permettent de supprimer w et w\_0

$$(a) \Rightarrow w = \sum_{m=1}^M \alpha_m y_m x_m \quad (\text{à vérifier})$$

$$(b) \Rightarrow \sum_{m=1}^M \alpha_m y_m = 0$$

En Substituant dans (2), le Lag ne depend plus que de alpha

$$\text{Lag}(\alpha) = \sum_{m=1}^M \alpha_m - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \alpha_i \alpha_j y_i y_j x_i^T x_j \quad (= 0)$$

à vérifier

Le pb s'écrit donc :

$$\alpha^* = \underset{\alpha_i | x_j}{\text{argmax}} \left\{ \sum_{m=1}^M \alpha_m - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \alpha_i \alpha_j y_i y_j x_i^T x_j \right\}$$

Sous les contraintes qqs m,  $\{\alpha_m > 0 \text{ et } \sum_{m=1}^M \alpha_m y_m = 0\}$

ON peut montrer facilement que la fonction à minimiser est une forme quadratique en alpha (TP)

Pb de programmation quadratique avec des contraintes + simples

Algorithmes efficaces et + spécifiques (gradient projeté, SMO, coordinate descent ...)

e1 nous dit que qqs m,  $\alpha_m(y_m v(x_m) - 1) = 0$

Autrement dit soit  $\alpha_m = 0$  soit  $y_m v(x_m) = 1$  : *caractéristique d'un vecteur support*)

Donc tous les  $\alpha_m$  liées aux individus qui ne sont pas des vecteurs supports sont nuls !

Peu de  $\alpha_m$  non nuls à calculer en pratique.

Une fois le pb dual résolu (par un algo itératif)  $\Rightarrow \alpha^*$

On calcule  $w^* = \sum_{m=1}^M \alpha_m y_m x_m$  (CL des vecteur supports)

et pour  $w_0$ , on choisit in individu m tq  $\alpha_m \neq 0$  (ie parmi les vecteurs supports) et résoud :

$$y_m(w_0^* + w^{*T} x_m) = 1$$

#### 4. Lien avec le coût de Hinge

Rappel:

$$e_{conv}^{hinge}(v(x), y) = \max(0, 1 - yv(x)) = \epsilon_m = \max(0, 1 - \mu) \Rightarrow SVM$$

$$R^{hinge} = \frac{1}{M} * \sum_{m=1}^M \max(0, 1 - y_m v(x_m)) = \sum_{m=1}^M \epsilon_m$$

Dans le cas linéairement séparable, toutes solution minimisant ce coût verifie donc qqs m:

$$1 - y_m v(x_m) \leq 0 \text{ ou } y_m v(x_m) \geq 1$$

Minimiser une version régularisée ridge du coût de Hinge revient donc à minimiser

$$\sum_{m=1}^M \epsilon_m + \alpha \|w\|^2$$

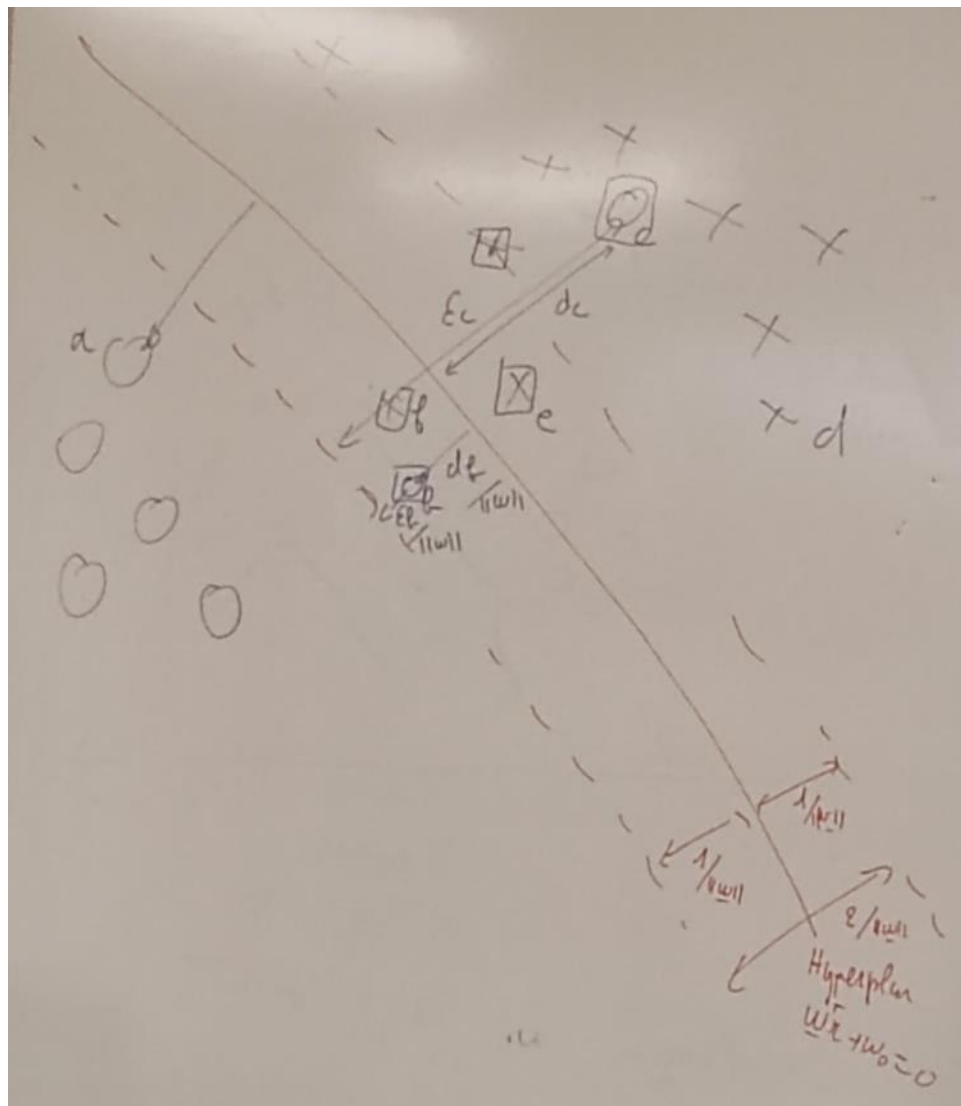
Sous contraintes  $y_m v(x_m) \geq 1$  : équivaut à la forme primal du pb SVM

### III) SVM LINEAIRE: CAS NON LINEAIREMENT SEPARABLE

#### 1. Objectif

Le classifieur linéaire reste pertinent mais ne peut pas classer correctement toutes les individus de la base d'apprentissage

- Fréquent en pratique
- On doit modifier l'approche précédente car la notion de marge dure n'est plus applicable ici.
- Pour cela, on se base sur le coût de Hinge quantifiant les erreurs de classification.



Que devient le coût de Hinge dans le cas non linéairement séparable pour un classifieur linéaire donné:

$$e_{conv}^{hinge}(v(x), y) = \max(0, 1 - yv(x)) = \epsilon_m = \max(0, 1 - \mu) \Rightarrow SVM$$

$$R^{hinge} = \frac{1}{M} * \sum_{m=1}^M \max(0, 1 - y_m v(x_m))$$

### Cas de figure:

Individu a :  $\frac{y_a v(x_a)}{\|w\|} \geq \frac{1}{\|w\|} \Rightarrow y_a v(x_a) \geq 1 \Rightarrow \epsilon_a = 0$  individu bien classé

b:  $0 < y_b v(x_b) < 1 \Rightarrow \epsilon_b = 1 - y_b v_b \Rightarrow 0 < \epsilon_b < 1$

1 bien classé mais dans la marge

c:  $y_c v(x_c) \leq -1 \Rightarrow \epsilon_c = 1 - y_c v_c \Rightarrow \epsilon_c > 2$  individu mal classé

d:  $y_d v(x_d) \geq 1 \Rightarrow \epsilon_d = 0$  bien classé

e:  $0 < y_e v(x_e) \leq -1 \Rightarrow \epsilon_e = 1 - y_e v_e \Rightarrow 0 < \epsilon_e < 1$

1 bien classé mais dans la marge

f:  $-1 < y_f v(x_f) \leq 0 \Rightarrow \epsilon_f = 1 - y_f v_f \Rightarrow 1 < \epsilon_f < 2$

2 mal classé mais dans la marge

$$R_{ridge}^{hinge} = \sum_{m=1}^M \epsilon_m + \lambda \|w\|^2$$

$$= C \sum_{m=1}^M \epsilon_m + \frac{1}{2} \|w\|^2 \quad C = \frac{1}{2\lambda} \text{ et le 2iem terme = régulation / taille marge}$$

le contrainte  $y_m v(x_m) \geq 1$  n'est plus vérifié pour tout m, elle doit être remplacée par une **contrainte plus souple**:  $y_m v(x_m) - \epsilon_m \geq 1$  et  $\epsilon_m \geq 0$

Le pb primal devient donc:

$$w, \theta_0, \epsilon^* = \operatorname{argmin} \left\{ C \sum_{m=1}^M \epsilon_m + \frac{1}{2} \|w\|^2 \right\} \text{ sous contrainte qqs } m, \epsilon_m \geq 1 -$$

$$y_m (w_0 + w^T x_m) \text{ et } \epsilon_m \geq 0$$

Par ailleurs la condition C1 devient :

$$\text{qqs } m, \alpha_m [y_m (w_0 + w^T x_m) - (1 - \epsilon_m)] = 0$$

- Soit  $\alpha_m = 0$

- Soit  $[y_m (w_0 + w^T x_m) - (1 - \epsilon_m)] = 0$

Les vecteurs supports sont ceux tq:

- Soit  $\epsilon = 0$  et  $y_m (w_0 + w^T x_m) = 1$  (sur la marge)

- Soit  $\epsilon \neq 0$  (Dans la marge ou mal classé)

C'est un hyperplan paramétré qui peut être déterminé par validation croisé ou à l'aide d'une base de test.

C petit  $\Rightarrow$  marge grand (forte tolérance aux erreurs d'apprentissage)

C grand  $\Rightarrow$  Contraire

## III) SVM NON LINEAIRE: ASTUCE NOYAU

### 1. Les classifieurs non linéaires

Les classifieurs linéaires forment une famille restreinte, pas toujours adaptés au pb.

On veut donc disposer de classifieurs plus généraux pour lesquels le séparateur n'est pas forcément un hyperplan mais une hypersurface d'équation  $v(x) = 0$  quelconque.

Parfois, une transformation évidente des variables explicatives suffit:

Ex: cas d'un séparateur cercle : visuellement un le séparateur peut être un cercle, donc on peut transformer l'ordonnancement des données avec un seuil r:

$$r = (X_1 - C)^2 + (X_2 - C)^2 \Rightarrow \text{un seuil sur } r \text{ suffit pour identifier}$$

C'est rarement édifiant en pratique,

=> Astuce du noyau

## 2. Astuce du noyau

Le principe est d'appliquer une transformation mathématique ( $\phi$ ) des données associant à chaque individu tel que les individus transformés soient linéairement séparables.

On aura alors plus cas ... ; classifieur linéaire des individus transformés:

$$v(x_m) = w^T \phi_m + w_0 \quad \phi_m = \phi(x_m)$$

La formulation duale des SVM permet de s'affranchir (en partie) du choix de  $\phi$ .

En effet, dans ce cas, après transformation, on a à résoudre:

$$\alpha^* = \underset{\alpha}{\operatorname{argmax}} \left\{ \sum_{m=1}^M \alpha_m - \frac{1}{2} \sum_{i=1}^M \sum_{j=1}^M \alpha_i \alpha_j y_i y_j \phi_i^T \phi_j \right\} \quad \phi_i^T \phi_j = <$$

Il suffit en pratique de connaître les valeurs des produits scalaires  $\phi_i^T \phi_j$  et pas forcément l'expression de  $\phi(x)$

On se contentera donc de choisir une fonction :  $k : R^M * R^M \rightarrow R$ , appelée **noyau**

Et qui devra simplement vérifier les propriétés d'un produit scalaire

$$k : R^M * R^M \rightarrow R \\ (x_i, x_j) \rightarrow K(x_i, x_j) = \phi_i^T \phi_j$$

Il existe des fonctions noyaux classiques:

- Noyau polynomial :  $K_\beta(x_i, x_j) = (x_i^T x_j + 1)^\beta$  degré  $\beta$
- Gaussien :  $K_\sigma(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma}}$
- Tan hyp :  $K_{h,c}(x_i, x_j) = \tanh(h x_i^T x_j + c)$
- Linéaire (SVM linéaire)  $K_\square(x_i, x_j) = x_i^T x_j$  ic pas de transformation  $\phi(x) = x \Rightarrow$  SVM linéaire

Un classifieur non linéaire correspond donc ici au choix d'une fonction noyau et à une classification linéaire par SVM

En pratique :

- On choisit  $K$
- On calcule les  $K(x_i, x_j)$
- On résout  $\alpha^* = \underset{\alpha}{\operatorname{argmax}} \dots$

**On ne peut pas remonter à  $w^*$  car  $w^* = \sum_{m=1}^M \alpha_m y_m \phi_m$   $\phi_m$  est inconnu**

Ce n'est pas un pb car on peut quand même faire l'étape de prédiction car elle dépend que des produits scalaires.

En effet:

$$f(x_{new}) = \operatorname{sign}(w_0^* + W^{*T} x_{new}) = \operatorname{sign} \left( \sum_{m=1}^M \alpha_m y_m \phi_m^T \phi(x_m + w_0^*) \right) \\ = \operatorname{sign} \left( \sum_{m=1}^M \alpha_m y_m K(x_m, x_{new}) + w_0^* \right)$$

Reste le  $w_0$

On choisit comme la solution de  $y_i v(\phi(x_i)) = 1$  où  $\phi(x_i)$  est un vecteur support

$$y_i (w^{*T} \phi(x_i) + w_0^*) = 1$$

$$y_i \left( \sum_{m=1}^M \alpha_m y_m K(x_m, x_i) + w_0^* \right) = 1$$