

# SVM

# Support Vector Machine

# Machine à Vecteur Support

Catherine ACHARD  
Institut des Systèmes Intelligents et de  
Robotique  
[catherine.achard@sorbonne-universite.fr](mailto:catherine.achard@sorbonne-universite.fr)

## Bibliographie

### Livres

C. Bishop, Pattern Recognition and Machine Learning, 2006

R. Duda, P. Hart et D. Stork, Pattern Classification, 2002

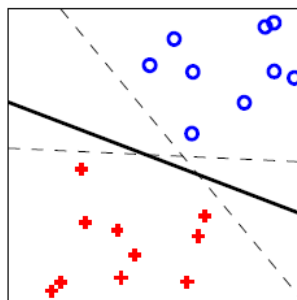
### Internet

<https://www.math.univ-toulouse.fr/~besse/Wikistat/pdf/st-m-app-svm.pdf>

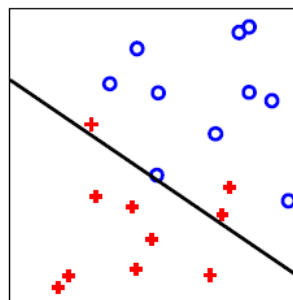
[https://www.lri.fr/~antoine/Courses/Master-ISI/ISI-10/Tr-cours-SVM\\_2014\\_2x2.pdf](https://www.lri.fr/~antoine/Courses/Master-ISI/ISI-10/Tr-cours-SVM_2014_2x2.pdf)

But : trouver une fonction de décision  $f(x)$  pour un problème à deux classes

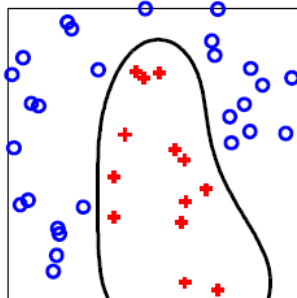
Image issue de Wikistat/pdf/st-m-app-svm.pdf



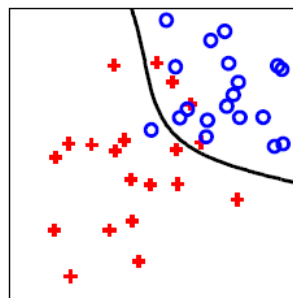
(a)



(b)



(c)



(d)

Exemples d'apprentissage :

$\{(x_i), i = 1, N\}$

Avec  $x_i \in \mathbb{R}^n$

Classes :

$\{(y_i), i = 1, N\}$

Avec  $y_i \in \{+1, -1\}$

Deux cas:

- Linéairement séparable (a et b)
- Non linéairement séparable (c et d)

Les données sont linéairement séparables s'il existe une fonction de décision  $f(x)$  linéaire (un hyperplan dans l'espace de dimension  $n$ ) permettant de les classer

$$f(x) = \mathbf{w}^T \cdot \mathbf{x} + b$$

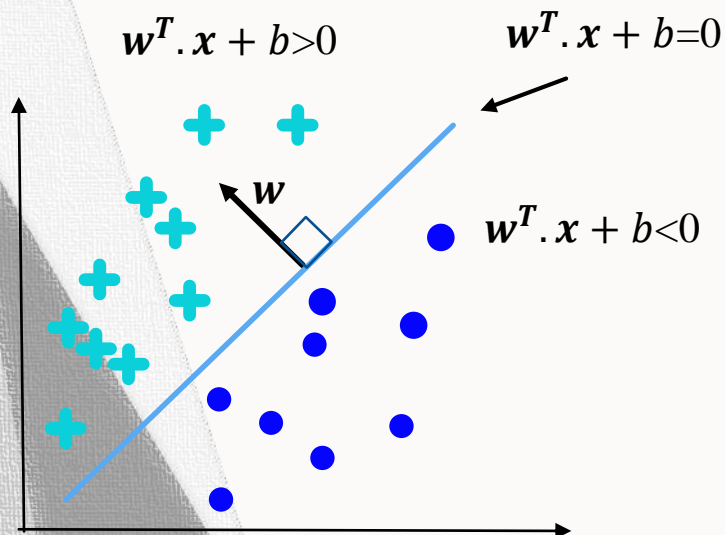
L'hyperplan séparateur vérifie

$$\mathbf{w}^T \cdot \mathbf{x}_i + b > 0 \text{ si } y_i = 1$$

$$\mathbf{w}^T \cdot \mathbf{x}_i + b < 0 \text{ si } y_i = -1$$

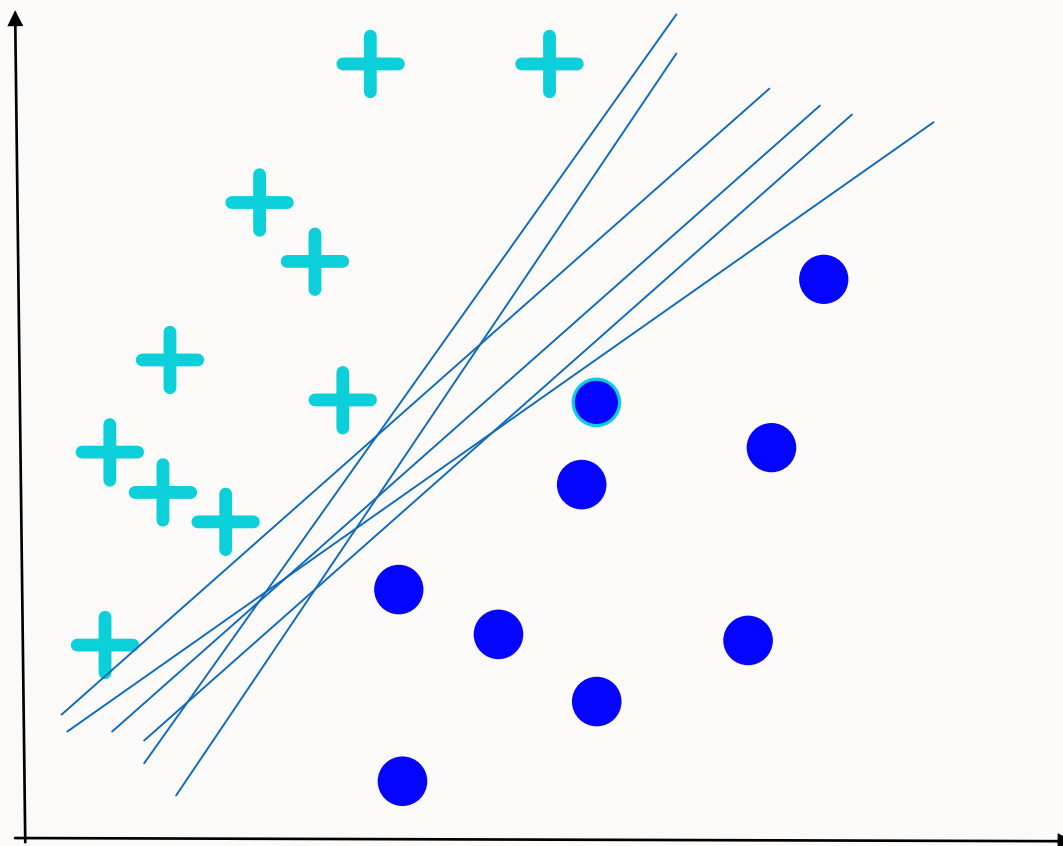
Ce qui revient à :

$$\forall i = 1, N \quad y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) > 0$$



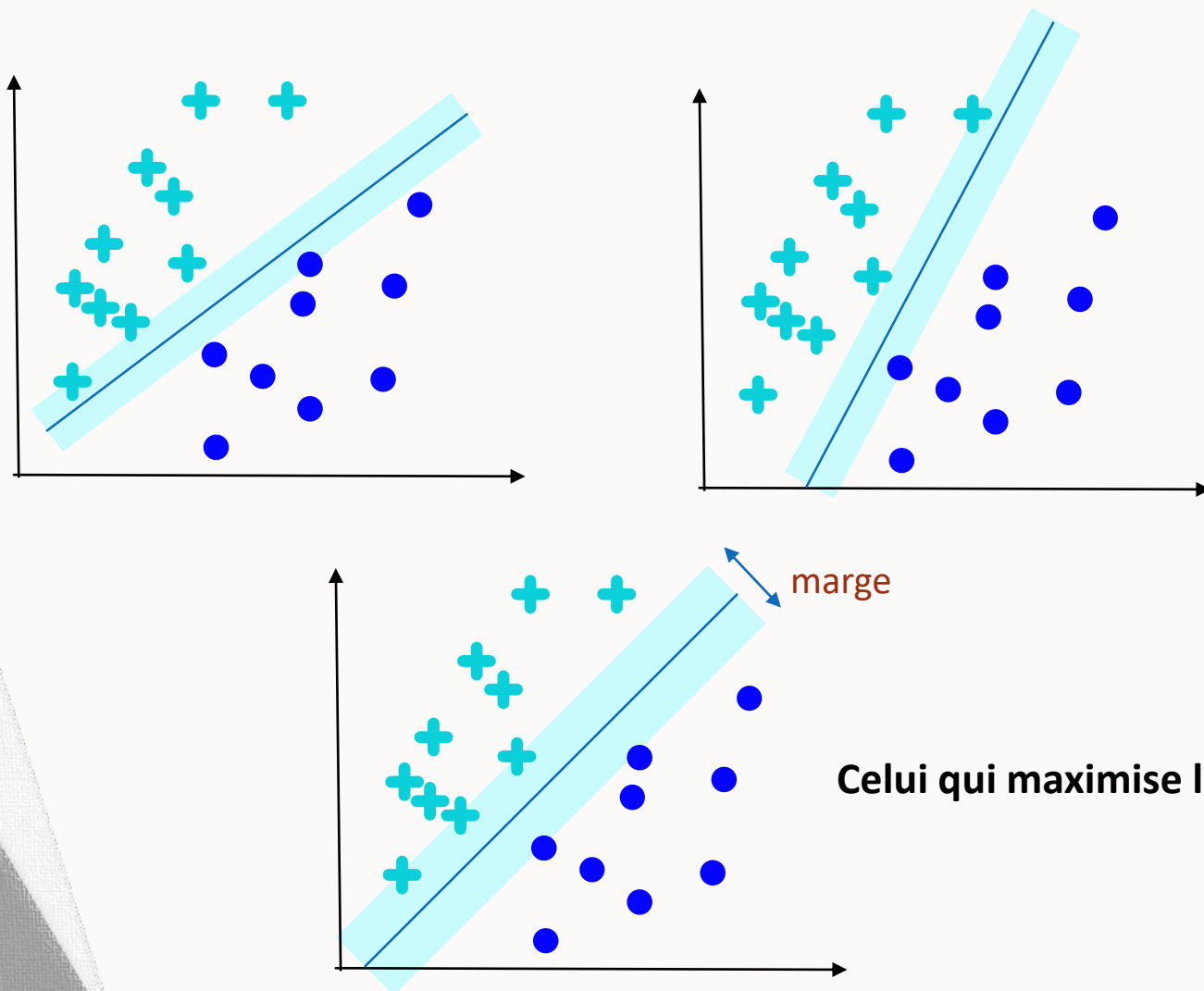
- $\mathbf{w}$  est un vecteur orthogonal à l'hyperplan
- $b$  est une constante

Il existe une infinité d'hyperplan, lequel choisir ?





Il existe une infinité d'hyperplan, lequel choisir ?



Toutes les équations  $k\mathbf{w}^T \cdot \mathbf{x}_i + kb = 0$  où  $k$  est une constante correspondent au même hyperplan.

On diminue l'espace des possibles en choisissant  $\mathbf{w}$  et  $b$  tq :

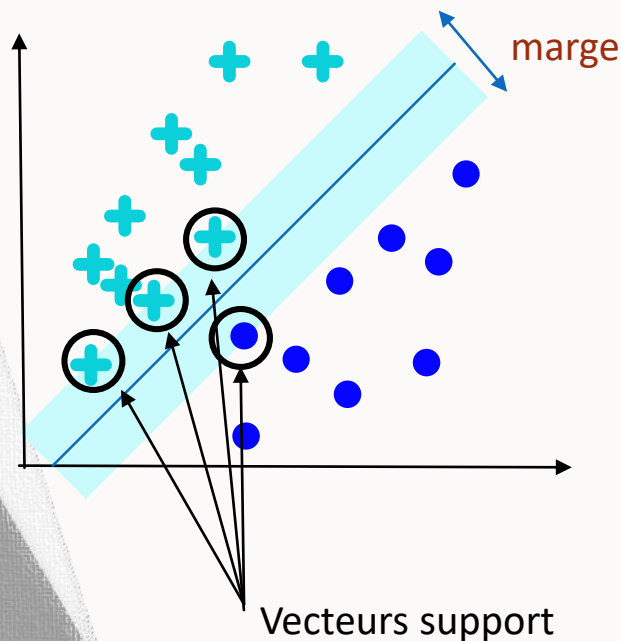
$$\min_{i=1,N} y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$$

Les exemples tq

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$$

sont appelés **vecteurs support**

Mais comment calculer la marge  $m$  ?

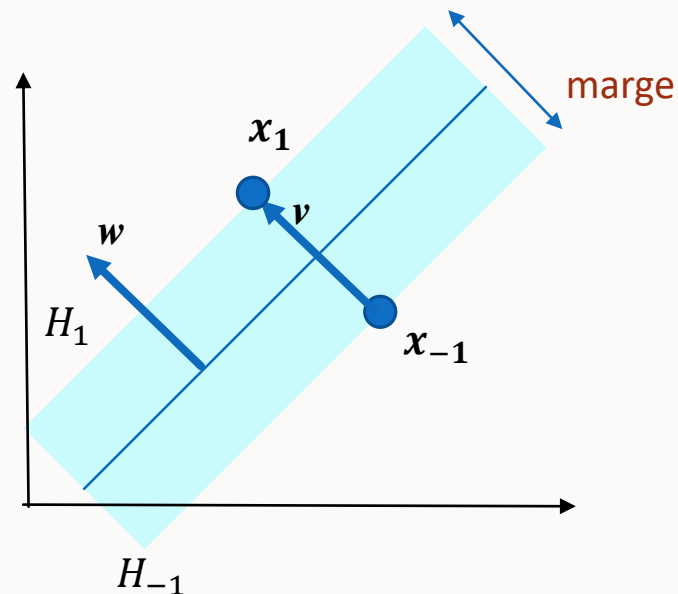


Comment trouver la marge  $m$  ?

Considérons les points  $x_{-1}$  et  $x_1$ , placés sur les hyperplans  $H_{-1}$  et  $H_1$  tq

$$\begin{aligned}w^T x_1 + b &= 1 \text{ (hyperplan } H_1) \\w^T x_{-1} + b &= -1 \text{ (hyperplan } H_{-1})\end{aligned}$$

$w$  est perpendiculaire à l'hyperplan. Soit  $v$  le vecteur, colinéaire à  $w$ , de norme  $m$  :  $v = \frac{w}{\|w\|}m$

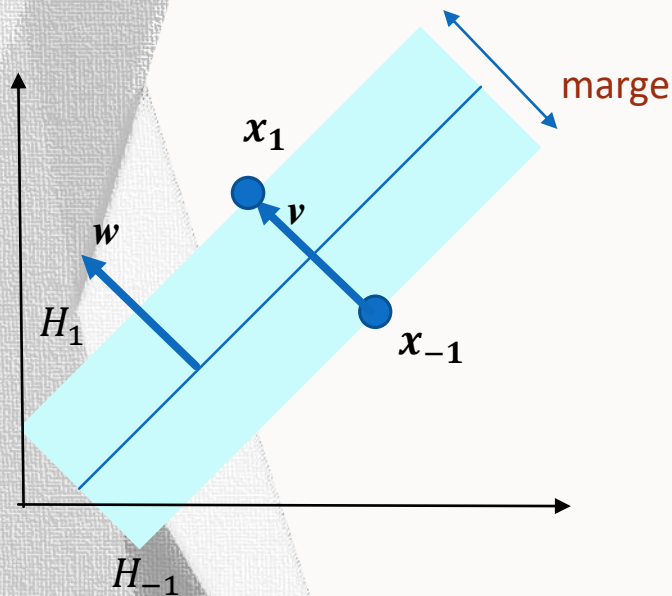




Considérons les points  $x_{-1}$  et  $x_1$ , placés sur les hyperplans  $H_{-1}$  et  $H_1$  tq

$$w^T x_1 + b = 1 \text{ et } w^T x_{-1} + b = -1$$

$w$  est perpendiculaire à l'hyperplan. Soit  $v$  le vecteur, colinéaire à  $w$ , de norme  $m$  :  $v = \frac{w}{\|w\|} m$



On a alors  $x_1 = x_{-1} + v$  avec  $v = \frac{w}{\|w\|} m$

$$w^T x_1 = w^T x_{-1} + \frac{w^T w}{\|w\|} m$$

$$w^T x_1 + b = w^T x_{-1} + b + \frac{w^T w}{\|w\|} m$$

$$1 = -1 + \|w\| m$$

$$m \|w\| = 2$$

Et donc la marge est :

$$m = \frac{2}{\|w\|}$$

Trouver l'hyperplan optimal revient donc au problème d'optimisation :

$$\max_{w,b} \frac{2}{\|w\|}$$

Sous contraintes :  $\forall i = 1, N \quad y_i (w^T \cdot x_i + b) \geq 1$

Ou encore,

$$\min_{w,b} \frac{1}{2} \|w\|^2 = \min_{w,b} \frac{1}{2} w^T w \quad \text{sc} \quad \forall i = 1, N \quad y_i (w^T \cdot x_i + b) \geq 1$$

Ce problème d'optimisation sous contrainte est un problème quadratique de la forme :

$$\min_z \frac{1}{2} z^T P z + q^T z \quad \text{sc} \quad G z \leq h$$

Il existe de nombreux solvers pour le résoudre.

**→ On arrive donc aux valeurs  $(w,b)$  de l'hyperplan optimal.**

Pour classer un nouvel exemple  $x$ , on utilise la **fonction de décision** :

$$f(x) = w^T \cdot x + b \quad \begin{cases} > 0 \Rightarrow y = 1 \\ < 0 \Rightarrow y = -1 \end{cases}$$

Résoudre

$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$
$$sc \quad \forall i = 1, N \quad y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1$$

Équivaut à résoudre

$$\min_z \frac{1}{2} \mathbf{z}^T \mathbf{P} \mathbf{z} + \mathbf{q}^T \mathbf{z}$$
$$sc \quad \mathbf{G} \mathbf{z} \leq \mathbf{h}$$

Trouver l'expression analytique de  $\mathbf{z}$ ,  $\mathbf{P}$ ,  $\mathbf{G}$ ,  $\mathbf{q}$  et  $\mathbf{h}$

Résoudre

$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$sc \quad \forall i = 1, N \quad y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1$$

Équivaut à résoudre

$$\min_z \frac{1}{2} \mathbf{z}^T \mathbf{P} \mathbf{z} + \mathbf{q}^T \mathbf{z}$$

$$sc \quad \mathbf{G} \mathbf{z} \leq \mathbf{h}$$

Trouver l'expression analytique de  $\mathbf{z}$ ,  $\mathbf{P}$ ,  $\mathbf{G}$ ,  $\mathbf{q}$  et  $\mathbf{h}$

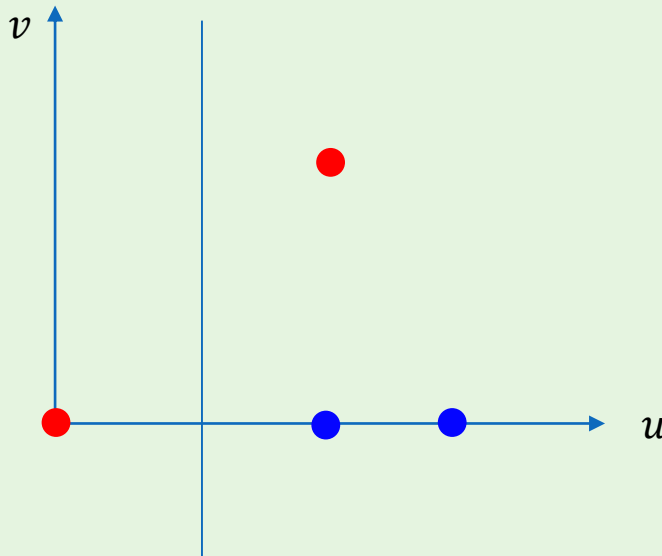
$$\mathbf{z} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^{n+1}, \mathbf{q} = \mathbf{0}_{n+1,1}, \mathbf{P} = \begin{bmatrix} 0 & \mathbf{0}_{1,n} \\ \mathbf{0}_{n,1} & \mathbf{I}_{n,n} \end{bmatrix}, \mathbf{G} = - \begin{bmatrix} y_1, y_1 \mathbf{x}_1^T \\ \vdots \\ y_N, y_N \mathbf{x}_N^T \end{bmatrix} \text{ et } \mathbf{h} = -\mathbf{1}_{N,1}$$

Considérons les points :

$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \text{ et } y = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$

1. Représentez ces points, sont-ils linéairement séparable ?
2. Représenter l'hyperplan de paramètre  $\mathbf{w}^T = [1,0]$  et  $b = -1$ . Est-il un bon classifieur pour ces données ?
3. Exprimer les contraintes de séparabilité pour chaque exemple
4. Calculer les paramètres  $(\mathbf{w}, b)$  de l'hyperplan optimal. Que vaut la marge ? Quels sont les vecteurs supports ?
5. Dans le cas où on souhaiterait résoudre le problème des SVM par un solveur de problème quadratique, déterminer les matrices  $P$ ,  $G$ ,  $q$ ,  $h$  à passer au solveur et la matrice  $z$  qu'il renverrait
6. Classer les points  $\begin{bmatrix} 0 \\ 1.5 \end{bmatrix}$  et  $\begin{bmatrix} 0 \\ 3.5 \end{bmatrix}$  en déterminant les valeurs de  $f(\mathbf{x})$ .





1. Représentez ces points, sont-ils linéairement séparable ?

Les points sont linéairement séparables

2. Représenter l'hyperplan de paramètre  $\mathbf{w}^T = [1, 0]$  et  $b = -1$ . Est-il un bon classifieur pour ces données ?

L'hyperplan a pour équation  $1 \cdot u + 0 \cdot v - 1 = 0$ , soit  $u = 1$ . Ce n'est pas un hyperplan séparateur

3. Exprimer les contraintes de séparabilité pour chaque exemple

Les contraintes de séparabilité sont

$$\forall i = 1, N \quad y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1$$

Et donc

$$\mathbf{x} = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \text{ et } \mathbf{y} = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix} \quad \begin{array}{ll} -b \geq 1 & (1) \\ -2w_1 - 2w_2 - b \geq 1 & (2) \\ 2w_1 + b \geq 1 & (3) \\ 3w_1 + b \geq 1 & (4) \end{array}$$

$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \text{ et } y = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix} \quad \begin{array}{l} -b \geq 1 \quad (1) \\ -2w_1 - 2w_2 - b \geq 1 \quad (2) \\ 2w_1 + b \geq 1 \quad (3) \\ 3w_1 + b \geq 1 \quad (4) \end{array}$$

4. Calculer les paramètres  $(\mathbf{w}, b)$  de l'hyperplan optimal. Que vaut la marge ? Quels sont les vecteurs supports ?

En combinant (1) et (3), on a :  $w_1 \geq 1$

En combinant (2) et (3), on a :  $w_2 \leq -1$

Or, on recherche  $\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} = \min_{\mathbf{w}, b} \frac{1}{2} (w_1^2 + w_2^2) \rightarrow w_1 = 1 \text{ et } w_2 = -1$

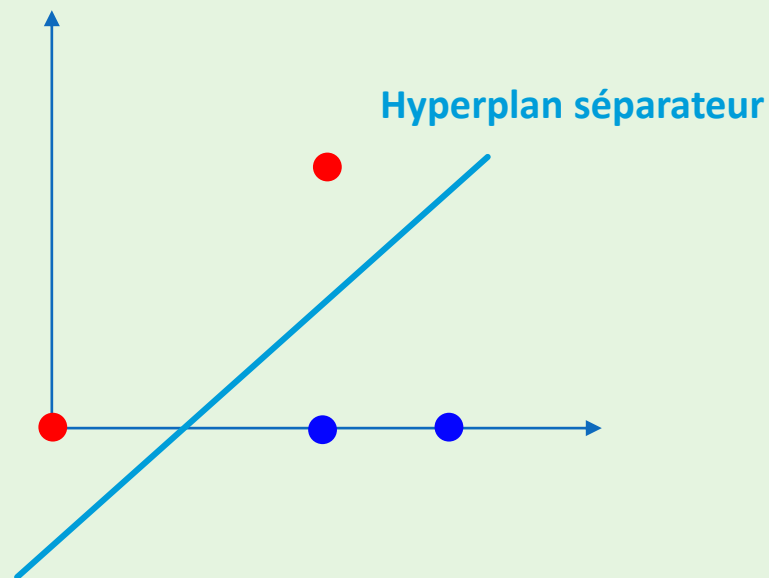
On obtient  $b$  en combinant (1) et (3) :

(1)  $\rightarrow b \leq -1$

(3)  $\rightarrow b \geq -1$



L'hyperplan est donc  
 $\mathbf{w}^T = [1, -1]$  et  $b = -1$



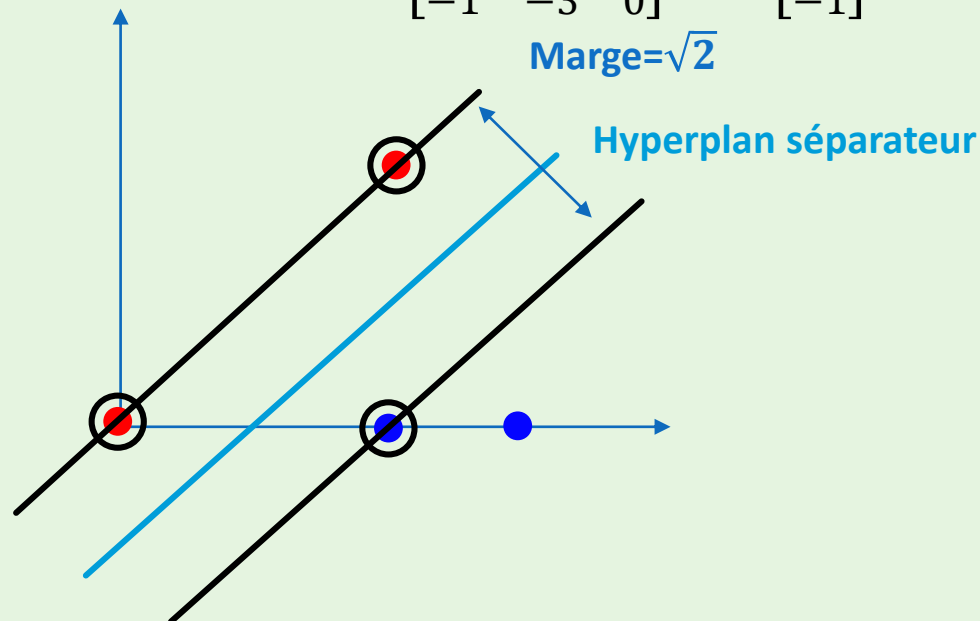
La marge vaut  $\frac{2}{\|w\|} = \frac{2}{\sqrt{2}} = \sqrt{2}$

Les vecteurs supports sont tq  $y_i(w^T \cdot x_i + b) = 1$  soit les points 1, 2 et 3

5. Dans le cas où on souhaiterait résoudre le problème des SVM par un solveur de problème quadratique, déterminer les matrices  $P$ ,  $G$ ,  $q$ ,  $h$  à passer au solveur et la matrice  $z$  qu'il renverrait

Pour la résolution avec un solveur de problème quadratique, on aurait

$$q = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, P = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, G = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 2 & 2 \\ -1 & -2 & 0 \\ -1 & -3 & 0 \end{bmatrix}, h = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \end{bmatrix} \text{ et } z = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix}$$

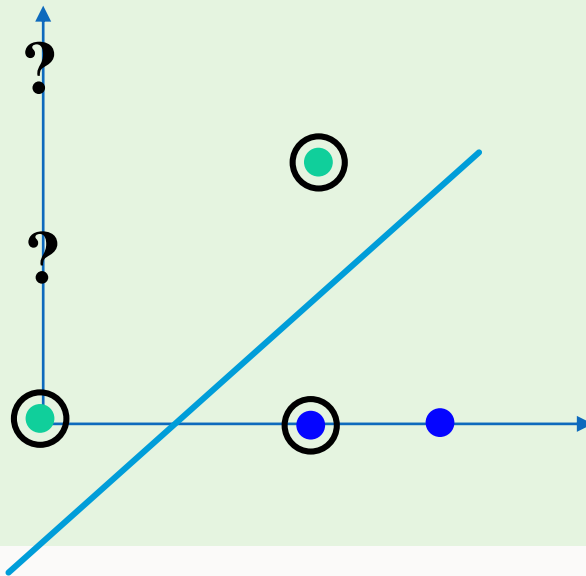


6. Classer les points  $\begin{bmatrix} 0 \\ 1.5 \end{bmatrix}$  et  $\begin{bmatrix} 0 \\ 3.5 \end{bmatrix}$  en déterminant les valeurs de  $f(\mathbf{x})$ .

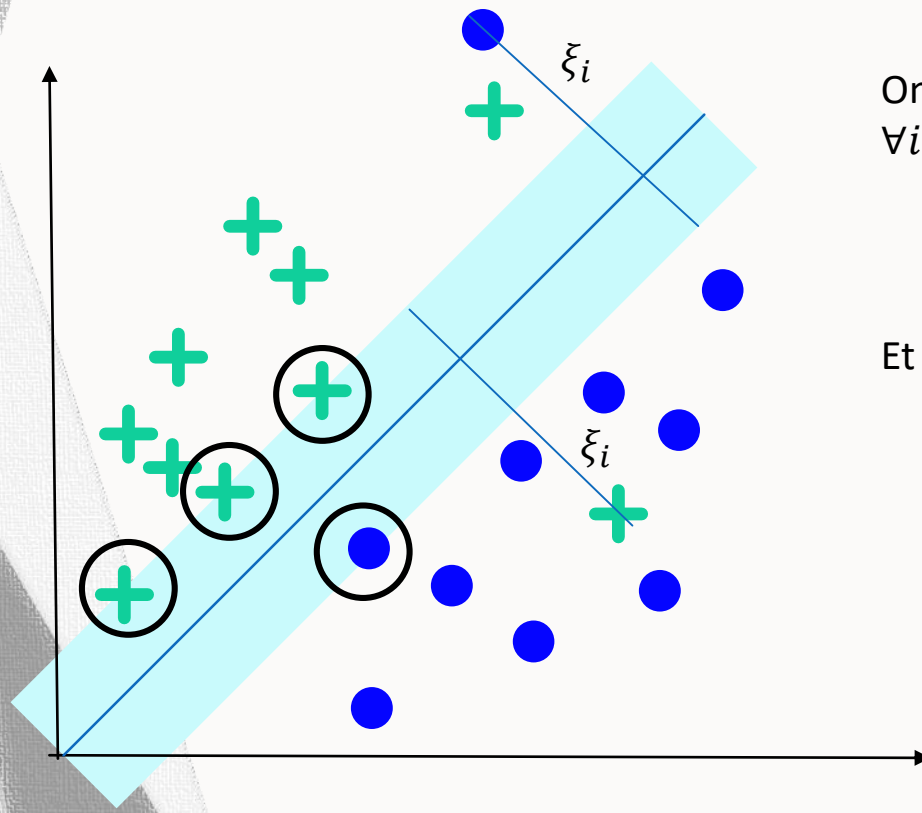
$$f(\mathbf{x}_1) = \mathbf{w}^T \cdot \mathbf{x}_1 + b = -2.5$$

$$f(\mathbf{x}_2) = \mathbf{w}^T \cdot \mathbf{x}_2 + b = -4.5$$

Les deux points appartiennent donc à la classe -1



Il se peut que les données ne soient pas linéairement séparables à cause de quelques exemples



On ajoute des variables d'ajustement  $\xi_i$  tq :  
 $\forall i = 1, N$

$$y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

et  $\xi_i \geq 0$

Et on souhaite que  $\sum_{i=1}^N \xi_i$  soit minimal



D'où la reformulation du problème :

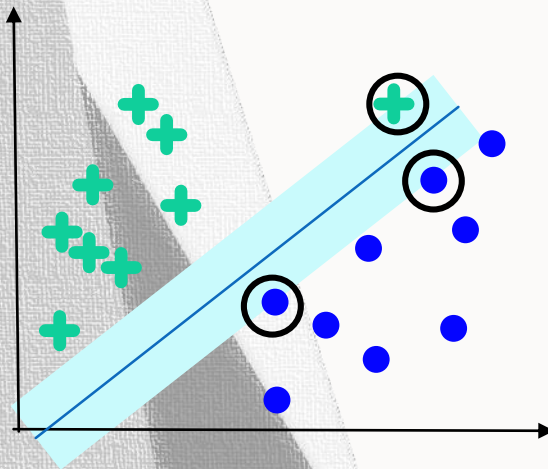
$$\min_{w,b,\xi_i} \left( \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \right)$$

sc  $\forall i = 1, N$

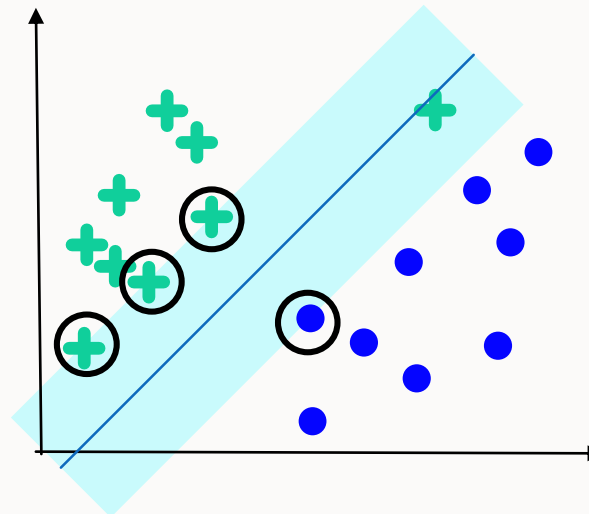
$$y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

$$\text{et } \xi_i \geq 0$$

- C règle le compromis entre une grande marge et le respect des contraintes
- Si C est grand, la marge sera petite pour que le maximum de point soit bien classé
- Si C est petit, la marge sera plus grande car on est plus laxiste sur les données mal classées
- C est à optimiser sur une base de validation



C grand



C petit

D'où la reformulation du problème :

$$\min_{w, b, \xi_i} \left( \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \right)$$

sc  $\forall i = 1, N$

$$y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

$$\text{et } \xi_i \geq 0$$

De la même manière, le problème se met sous une forme quadratique :

$$\min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T \mathbf{P} \mathbf{z} + \mathbf{q}^T \mathbf{z}$$

$$\text{sc } \mathbf{G} \mathbf{z} \leq \mathbf{h}$$

Exercice :

Déterminer les expressions analytiques de  $\mathbf{z}$ ,  $\mathbf{P}$ ,  $\mathbf{G}$ ,  $\mathbf{q}$  et  $\mathbf{h}$

D'où la reformulation du problème :

$$\min_{w, b, \xi_i} \left( \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \right)$$

$$sc \quad \forall i = 1, N \quad y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1 - \xi_i \text{ et } \xi_i \geq 0$$

De la même manière, le problème se met sous une forme quadratique :

$$\min_{\mathbf{z}} \frac{1}{2} \mathbf{z}^T \mathbf{P} \mathbf{z} + \mathbf{q}^T \mathbf{z} \quad sc \quad \mathbf{G} \mathbf{z} \leq \mathbf{h}$$

Et, en identifiant,

$$\mathbf{z} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \\ \xi_1 \\ \vdots \\ \xi_N \end{bmatrix}, \mathbf{P} = \begin{bmatrix} 0 & \mathbf{0}_{1,n} & \mathbf{0}_{1,N} \\ \mathbf{0}_{n,1} & \mathbf{I}_{n,n} & \mathbf{0}_{n,N} \\ \mathbf{0}_{N,1} & \mathbf{0}_{N,n} & \mathbf{0}_{N,N} \end{bmatrix}, \mathbf{q} = \begin{bmatrix} \mathbf{0}_n^T \\ C \cdot \mathbf{1}_{N,1}^T \end{bmatrix},$$

$$\mathbf{G} = \begin{bmatrix} -y_1, -y_1 \mathbf{x}_1^T & & \\ & \vdots & \\ -y_N, -y_N \mathbf{x}_N^T & & -\mathbf{I}_{N,N} \\ & \mathbf{0}_{N,n+1} & -\mathbf{I}_{N,N} \end{bmatrix}, \mathbf{h} = \begin{bmatrix} -\mathbf{1}_{N,1} \\ \mathbf{0}_{N,1} \end{bmatrix}$$

En appliquant ces équations au problème précédant avec  $C=1$ , on trouve :

$$z = \begin{bmatrix} -1.0000 \\ 0.9999 \\ -0.9999 \\ 0.0000 \\ 0.0000 \\ 0.0003 \\ 0.0003 \end{bmatrix}$$

Que peut-on conclure ?

En appliquant ces équations au problème précédant avec  $C=1$ , on trouve :

$$z = \begin{bmatrix} -1.0000 \\ 0.9999 \\ -0.9999 \\ 0.0000 \\ 0.0000 \\ 0.0003 \\ 0.0003 \end{bmatrix}$$

Que peut-on conclure ?

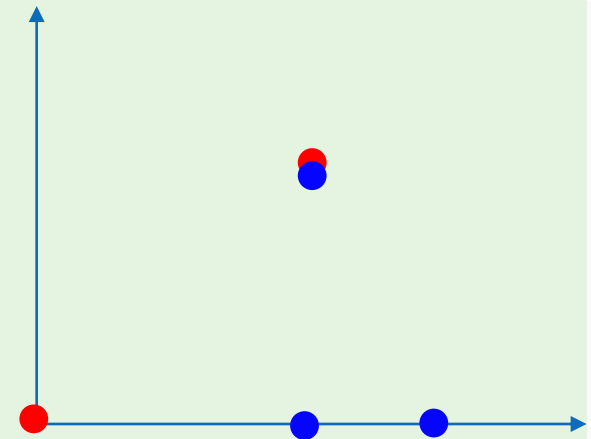
On retrouve presque le même hyperplan :  $\mathbf{w}^T = [1, -1]$  et  $b = -1$

Tous les  $\xi_i$  sont très petits  $\rightarrow$  les exemples respectent bien la contrainte



On ajoute maintenant un point aux données :

$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \\ 2 & 1.9 \end{bmatrix} \text{ et } y = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$



En faisant tourner le programme, on obtient :

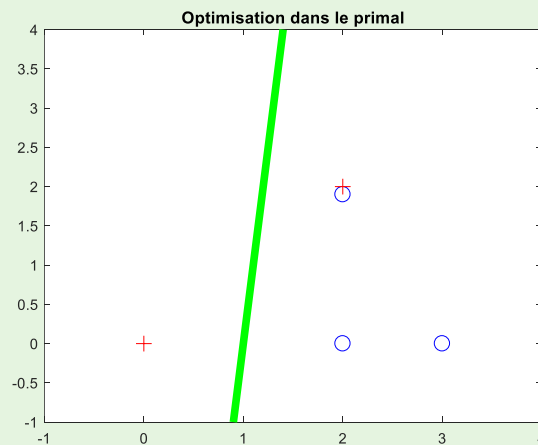
$$\text{Avec } C=1, z = \begin{bmatrix} -1.0000 \\ 1 \\ -0.1 \\ 0.0000 \\ 1.8000 \\ 0.0000 \\ 0.1900 \end{bmatrix} \text{ et avec } C=1000, z = \begin{bmatrix} -1.0000 \\ 20 \\ -20 \\ 0.0000 \\ 0.0000 \\ 0.0000 \\ 0.0000 \end{bmatrix}$$

Commentez

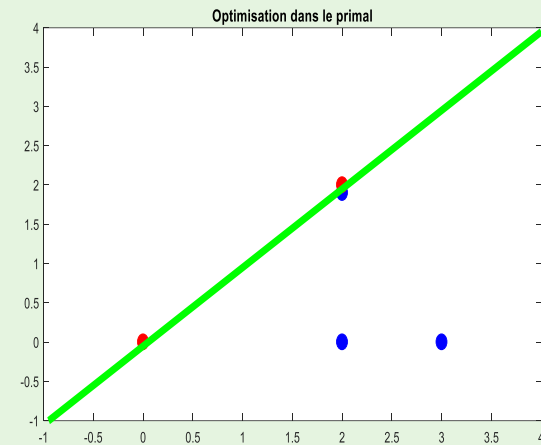
$C=1 \rightarrow$  On donne moins de poids au respect des contraintes  $\rightarrow$  La marge est plus grande mais le second point viole la contrainte

$C=1000 \rightarrow$  on donne un très fort poids au respect des contraintes  $\rightarrow$  Elles sont respectées mais la marge est très petite

$C=1$



$C=1000$



Nous avons vu que l'expression primale des SVM est :

$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$sc \quad \forall i = 1, N \quad y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) \geq 1$$

On peut résoudre ce problème en utilisant le lagrangien :

$$L(\mathbf{w}, b, \alpha) \quad \begin{cases} = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) - 1] \\ = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)] \end{cases}$$

Les  $\alpha_i$  sont appelés **multiplicateurs de Lagrange** et  $\forall i, \alpha_i \geq 0$ .

On recherche :

$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \max_{\alpha_i} \sum_{i=1}^N \alpha_i [1 - y_i(\mathbf{w}^T \cdot \mathbf{x}_i + b)]$$

**Pourquoi ?**

**Pour des données linéairement séparables :**

- Si l'hyperplan est tq la **contrainte n'est pas respectée en  $x_i$** 
  - $[1 - y_i(\mathbf{w}^T \cdot \mathbf{x}_i + b)] > 0$
  - Le  $\alpha_i$  correspondant devient très grand pour maximiser  $\alpha_i [1 - y_i(\mathbf{w}^T \cdot \mathbf{x}_i + b)]$
  - $\mathbf{w}, b$  vont évoluer de manière à minimiser le Lagrangien et vérifier la contrainte
- Si la contrainte est **strictement respectée en  $x_i$** 
  - $[1 - y_i(\mathbf{w}^T \cdot \mathbf{x}_i + b)] < 0$
  - $\alpha_i = 0$  pour maximiser  $\alpha_i [1 - y_i(\mathbf{w}^T \cdot \mathbf{x}_i + b)]$
- Si  $\mathbf{x}_i$  est un vecteur support
  - $1 - y_i(\mathbf{w}^T \cdot \mathbf{x}_i + b) = 0$
  - $\alpha_i > 0$

Pour conclure, pour un cas linéairement séparable,  
 $\alpha_i = 0$  → L'exemple  $\mathbf{x}_i$  vérifie la contrainte  
 $\alpha_i > 0$  → vecteur support

D'autre part, résoudre:

$$\min_{w,b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \max_{\alpha_i} \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)]$$

Revient à résoudre

$$\max_{\alpha_i} \min_{w,b} \underbrace{\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)]}_{L(\mathbf{w}, b, \alpha)}$$

$$\max_{\alpha_i} \min_{w,b} L(\mathbf{w}, b, \alpha)$$



$$\max_{\alpha_i} \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)]$$

En dérivant  $L(\mathbf{w}, b, \alpha)$  par rapport à  $\mathbf{w}$  et  $b$ ,

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = 0 \rightarrow ?$$

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = 0 \rightarrow ?$$

Rappel sur les matrices

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$$

$$\frac{\partial a^T b}{\partial a} = \frac{\partial b^T a}{\partial a} = b$$

$$\frac{\partial a^T \mathbf{B} a}{\partial a} = 2 \mathbf{B} a$$

$$\max_{\alpha_i} \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)]$$

En dérivant  $L(\mathbf{w}, b, \alpha)$  par rapport à  $\mathbf{w}$  et  $b$ ,

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad (1)$$

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = 0 \rightarrow \sum_{i=1}^N \alpha_i y_i = 0 \quad (2)$$

Rappel sur les matrices

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$$

$$\frac{\partial a^T b}{\partial a} = \frac{\partial b^T a}{\partial a} = b$$

$$\frac{\partial a^T \mathbf{B} a}{\partial a} = 2 \mathbf{B} a$$

En dérivant  $L(\mathbf{w}, b, \alpha)$  par rapport à  $\mathbf{w}$  et  $b$ ,

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad (1)$$

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = 0 \rightarrow \sum_{i=1}^N \alpha_i y_i = 0 \quad (2)$$

Or,

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)]$$

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i y_i \mathbf{w}^T \cdot \mathbf{x}_i - \sum_{i=1}^N \alpha_i y_i b$$

On remplace  $\mathbf{w}$   
par sa valeur  
(eq 1)

On remplace  $\mathbf{w}$   
par sa valeur  
(eq 1)

=0 (eq 2)

Rappel sur les matrices

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$$

$$\frac{\partial a^T b}{\partial a} = \frac{\partial b^T a}{\partial a} = b$$

$$\frac{\partial a^T \mathbf{B} a}{\partial a} = 2 \mathbf{B} a$$

En dérivant  $L(\mathbf{w}, b, \alpha)$  par rapport à  $\mathbf{w}$  et  $b$ ,

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial \mathbf{w}} = 0 \rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \quad (1)$$

$$\frac{\partial L(\mathbf{w}, b, \alpha)}{\partial b} = 0 \rightarrow \sum_{i=1}^N \alpha_i y_i = 0 \quad (2)$$

Or, 
$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b)]$$

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \alpha_i y_i \mathbf{w}^T \cdot \mathbf{x}_i - \sum_{i=1}^N \alpha_i y_i b$$

Soit, en introduisant (1) et (2) dans le lagrangien,

$$L(\alpha) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^N \alpha_i - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$L(\alpha) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^N \alpha_i$$

On souhaite maximiser le lagrangien par rapport à  $\alpha$ ,

$$\max_{\alpha_i} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad \text{sc} \quad \sum_{i=1}^N \alpha_i y_i = 0 \text{ (eq 2) et } \alpha_i \geq 0$$

Rappel sur les matrices

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$(\mathbf{A} + \mathbf{B})^T = \mathbf{A}^T + \mathbf{B}^T$$

$$\frac{\partial a^T b}{\partial a} = \frac{\partial b^T a}{\partial a} = b$$

$$\frac{\partial a^T \mathbf{B} a}{\partial a} = 2 \mathbf{B} a$$

$$\max_{\alpha_i} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad \text{sc } \sum_{i=1}^N \alpha_i y_i = 0 \text{ et } \alpha_i \geq 0$$

revient à

$$\min_{\alpha_i} \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j - \sum_{i=1}^N \alpha_i \quad \text{sc } \sum_{i=1}^N \alpha_i y_i = 0 \text{ et } \alpha_i \geq 0$$

Ce problème peut aussi être mis sous une forme quadratique.

→ On trouve les valeurs de  $\alpha_i$

Mais alors, quel est l'intérêt ?

- L'astuce du noyau permet de résoudre des problèmes non linéaires



Concrètement, la résolution du problème quadratique nous amène aux valeurs de  $\alpha_i$

Pour les problèmes linéairement séparable, on en déduit les valeurs de  $\mathbf{w}$  avec :

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = \sum_{\alpha_i > 0} \alpha_i y_i \mathbf{x}_i$$

Puis la valeur de  $b$  :

**En se plaçant sur un des vecteurs supports** (un des points tq  $\alpha_i > 0$ )

Soit

$$y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) = 1$$

Ou

$$y_i y_i (\mathbf{w}^T \cdot \mathbf{x}_i + b) = y_i$$

Et donc,

$$(\mathbf{w}^T \cdot \mathbf{x}_i + b) = y_i \text{ car } y_i = \pm 1$$

$$\mathbf{b} = y_i - \mathbf{w}^T \cdot \mathbf{x}_i$$

$$\mathbf{b} = y_i - \left( \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^T \cdot \mathbf{x}_i$$

Pour classer un nouvel exemple  $\mathbf{x}$ , la fonction de décision est :

$$f(\mathbf{x}) = \text{signe}(\mathbf{w}^T \cdot \mathbf{x}_i + b)$$

Ou encore

$$f(\mathbf{x}) = \text{signe} \left( \sum_{\alpha_i > 0} \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_i + \mathbf{b} \right)$$

Reprenons l'exercice précédant :

$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \text{ et } y = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$

Le solver renvoie  $\alpha = \begin{bmatrix} 0.5 \\ 0.5 \\ 1 \\ 0 \end{bmatrix}$ . Quels sont les vecteurs support ? Retrouver les paramètres  $\mathbf{w}$  et  $\mathbf{b}$  de l'hyperplan

$$x = \begin{bmatrix} 0 & 0 \\ 2 & 2 \\ 2 & 0 \\ 3 & 0 \end{bmatrix} \text{ et } y = \begin{bmatrix} -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}$$

Reprenons l'exercice précédant. Le solver renvoie  $\alpha = \begin{bmatrix} 0.5 \\ 0.5 \\ 1 \\ 0 \end{bmatrix}$ . Quels sont les vecteurs support ? Donner l'équation de la fonction de décision

Les vecteurs supports correspondent aux positifs. Donc les points 1, 2 et 3.  
L'hyperplan séparateur est obtenu avec :

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\mathbf{b} = y_i - \mathbf{w}^T \cdot \mathbf{x}_i \text{ pour un des vecteurs supports}$$

$$\text{Soit, } \mathbf{w} = -0.5 \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.5 \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

En prenant le premier point  $\mathbf{b} = -1$

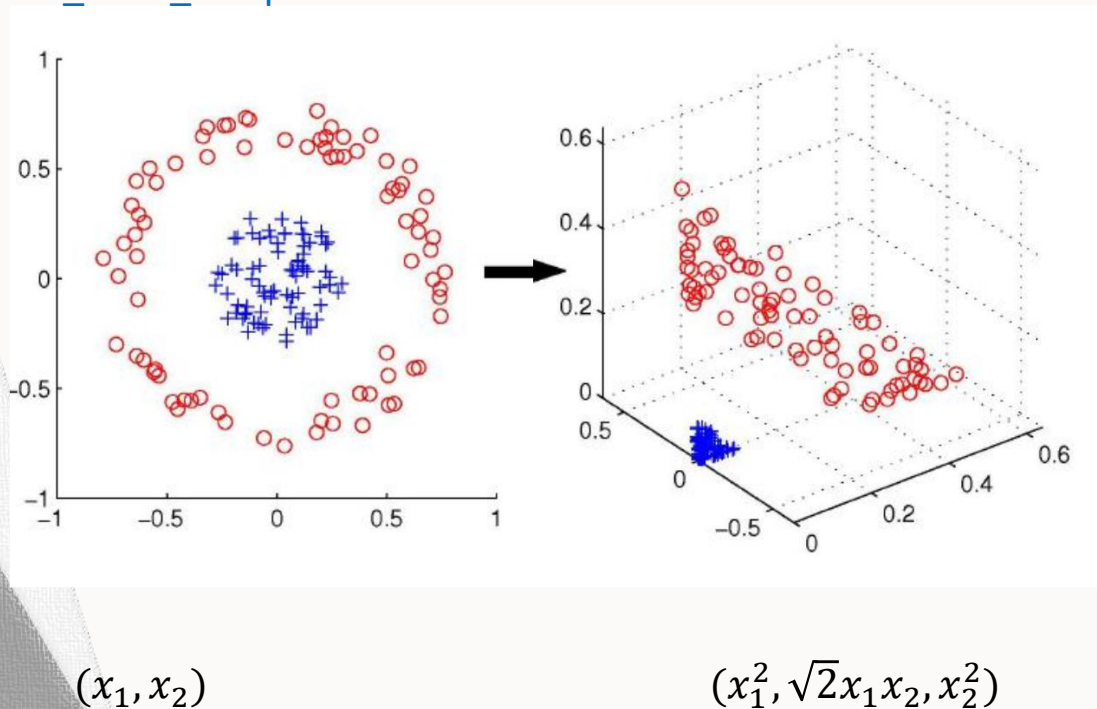
On retrouve bien les mêmes résultats que précédemment

Ce problème est mis sous une forme quadratique. Mais alors, quel est l'intérêt ?

- **L'astuce du noyau permet de résoudre des problèmes non linéaires**

Une solution est de transformer les données pour les rendre linéairement séparables, Ceci demande souvent d'augmenter la dimension du problème

Figure issue de [https://www.lri.fr/~antoine/Courses/Master-ISI/ISI-10/Tr-cours-SVM\\_2014\\_2x2.pdf](https://www.lri.fr/~antoine/Courses/Master-ISI/ISI-10/Tr-cours-SVM_2014_2x2.pdf)





Ainsi, les données  $\mathbf{x}$  seront transformées en  $\Phi(\mathbf{x})$ .

En reprenant la formulation duale :

$$\max_{\alpha_i} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$sc \sum_{i=1}^N \alpha_i \alpha y_i = 0 \text{ et } \alpha_i \geq 0$$

Dans le nouvel espace,

$$\max_{\alpha_i} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$$

$$sc \sum_{i=1}^N \alpha_i \alpha y_i = 0 \text{ et } \alpha_i \geq 0$$

$\Phi(\mathbf{x})$  peut être de grande dimension mais il n'a pas besoin d'être calculé dès lors que l'on peut calculer le produit scalaire  $\Phi(\mathbf{x}_i)^T \Phi(\mathbf{x}_j)$  de manière rapide avec une fonction noyau

Dans le cas précédant,

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \text{ et } \Phi(\mathbf{x}) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}$$

Or,

$$\Phi(\mathbf{x})^T \cdot \Phi(\mathbf{z}) = \begin{bmatrix} x_1^2 & \sqrt{2}x_1x_2 & x_2^2 \end{bmatrix} \cdot \begin{bmatrix} z_1^2 \\ \sqrt{2}z_1z_2 \\ z_2^2 \end{bmatrix}$$

$$\Phi(\mathbf{x})^T \cdot \Phi(\mathbf{z}) = x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2$$

$$\Phi(\mathbf{x})^T \cdot \Phi(\mathbf{z}) = (x_1 z_1 + x_2 z_2)^2$$

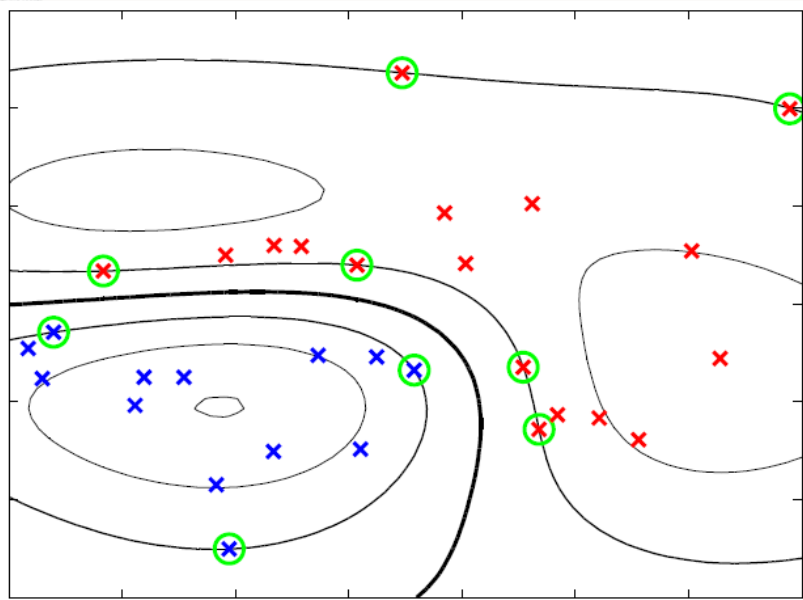
$$\Phi(\mathbf{x})^T \cdot \Phi(\mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$$

On a donc  $\Phi(\mathbf{x})^T \cdot \Phi(\mathbf{z}) = K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2$

➔ Le produit scalaire peut être calculé sans passer par la transformation des données

Les noyaux fréquemment utilisés sont

- Les noyaux polynomiaux d'ordre  $p$  :  $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + 1)^p$
- Les noyaux gaussiens :  $K(\mathbf{x}, \mathbf{z}) = e^{-\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2}}$

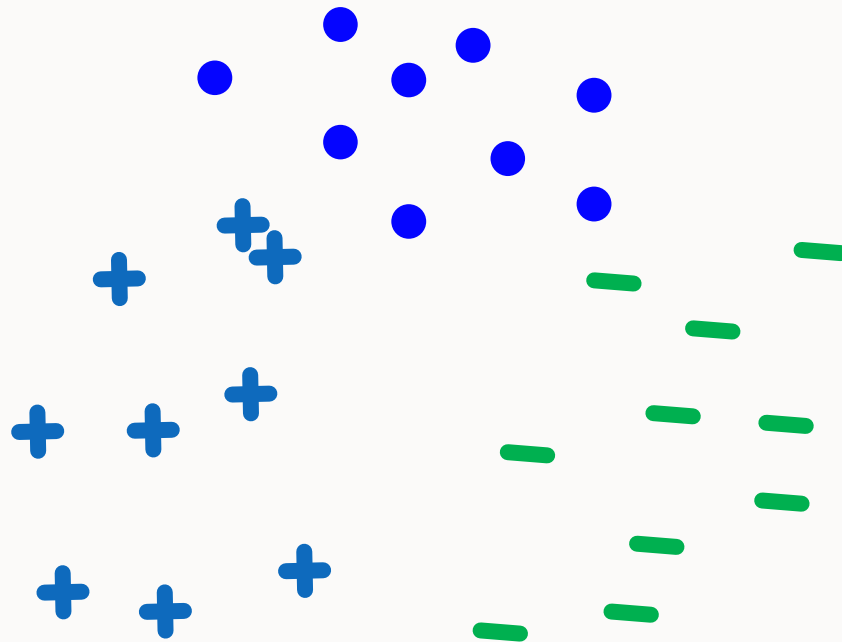


Séparation obtenue avec un noyau gaussien  
Les points verts représentent les vecteurs support

Représentation dans l'espace original

Figure issue de C. Bishop, Pattern Recognition and Machine Learning, 2006

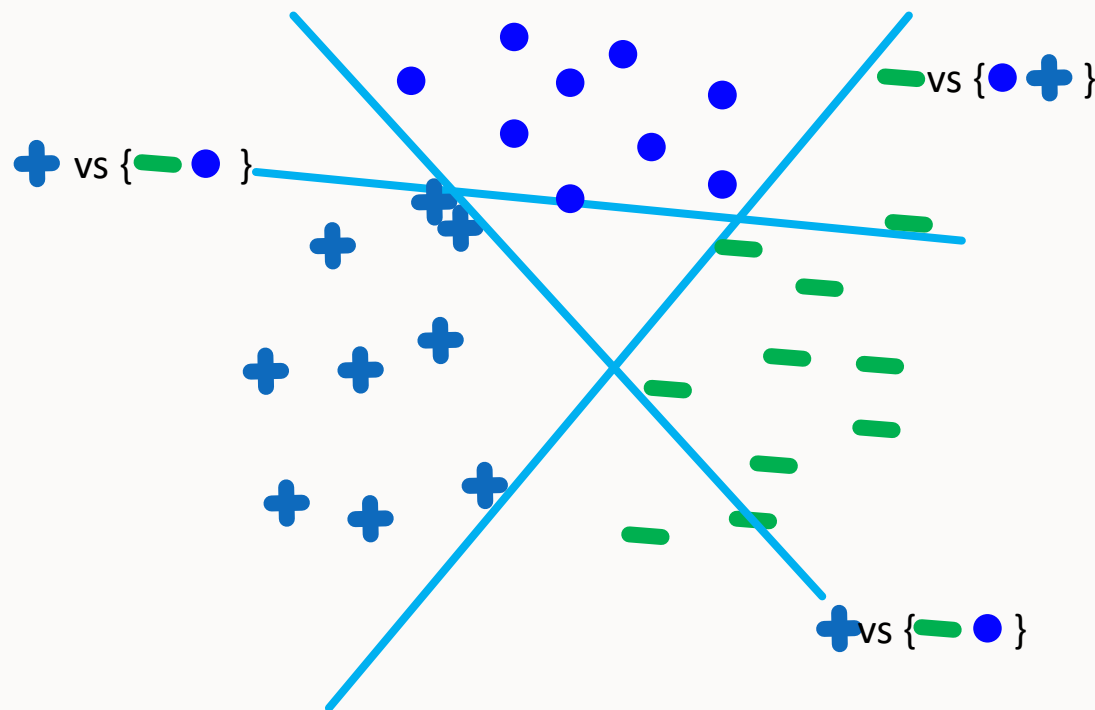
Les SVM ont été conçus pour résoudre des problèmes binaires.  
Comment faire pour les problèmes multi-classes ?



## L'approche 1 contre tous

Dans un problème à K classes, K SVM qui discriminent une classe contre les autres

Exemple



On entraîne 3 SVM binaires

+ vs { - ● }

- vs { ● + }

+ vs { - ● }

En test, l'étiquette est donnée par le SVM avec la plus grande fonction de décision  $f(\mathbf{x}) > 0$  (l'exemple classé  $> 0$  qui répond le mieux)

## Inconvénients de l'approche 1 contre tous

- Pour chaque SVM, il y a beaucoup plus d'exemples négatifs



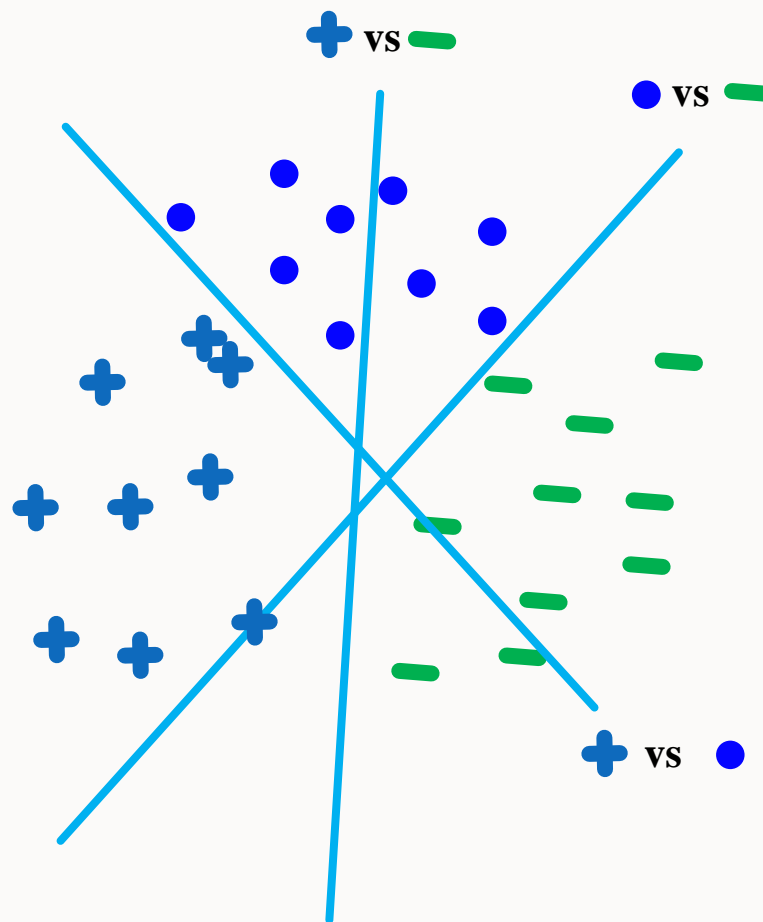
## L'approche 1 contre 1

On construit un SVM pour chaque paire d'étiquettes  $\rightarrow K(K-1)/2$  SVM

Exemple

- vs —
- + vs ●
- + vs —

La prédiction correspond au label qui a gagné le plus de « duels »



### **Inconvénients de l'approche 1 contre 1**

- Très long :  $K(K-1)/2$  SVM à apprendre

### **Avantage de l'approche 1 contre 1**

- Moins de déséquilibre entre classes