

Mathématiques de Machine Learning 12 MAPS 3  
2024-2025

Exercice n° 3:

1) Il s'agit d'un problème de classification.

$$2) \hat{f} \in \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{f(x_i) \neq y_i}$$

3)  $\forall x \in \mathcal{X}, c \in \{-1, 1\}$   $f_x(x)$

$$P(X=x, Y=c) = P(X=x) P(Y=c | X=x)$$

de plus  $P(X=x, Y=c) = \underbrace{P(Y=c)}_{\pi_c} \underbrace{P(X=x | Y=c)}_{f_{X|Y=c}(x)}$

donc  $\forall x, \forall c, \underbrace{f_x(x)}_{>0} P(Y=c | X=x) = \pi_c f_{X|Y=c}(x)$

d'où:  $P(Y=c | X=x) = \frac{\pi_c f_{X|Y=c}(x)}{f_x(x)}$



or, d'après la formule des probabilités totales,

$$f_X(x) = \pi_1 f_{X|Y=1}(x) + \pi_{-1} f_{X|Y=-1}(x)$$

$$\text{d'où, } \forall x \in \mathbb{R} \quad P(Y=c|X=x) = \frac{\pi_c f_{X|Y=c}(x)}{\pi_1 f_{X|Y=1}(x) + \pi_{-1} f_{X|Y=-1}(x)}$$

$$4) \underline{P_{\hat{\pi}_1}}(y_1, \dots, y_n) = \prod_{i=1}^n \hat{\pi}_1^{1_{y_i=1}} (1 - \hat{\pi}_1)^{1 - 1_{y_i=1}}$$

donc en prenant le log :

$$\log P_{\hat{\pi}_1}(y_1, \dots, y_n) = \sum_{i=1}^n \left( \frac{1_{y_i=1}}{\hat{\pi}_1} \log(\hat{\pi}_1) + \left(1 - \frac{1_{y_i=1}}{\hat{\pi}_1}\right) \log(1 - \hat{\pi}_1) \right)$$

$$\frac{\partial}{\partial \hat{\pi}_1} \log P_{\hat{\pi}_1}(y_1, \dots, y_n) = \sum_{i=1}^n \left( \frac{1_{y_i=1}}{\hat{\pi}_1} - \frac{1 - 1_{y_i=1}}{1 - \hat{\pi}_1} \right)$$

$$\frac{\partial}{\partial \hat{\pi}_1} \log P_{\hat{\pi}_1}(y_1, \dots, y_n) = 0 \quad \text{iff}$$

$$\sum_{i=1}^n \left( \frac{1_{y_i=1}}{\hat{\pi}_1} - \frac{1 - 1_{y_i=1}}{1 - \hat{\pi}_1} \right) = 0$$

$$\text{iff } \frac{1 - \hat{\pi}_1}{\hat{\pi}_1} \left( \sum_{i=1}^n 1_{y_i=1} \right) = n - \sum_{i=1}^n 1_{y_i=1}$$



ipp

$$\sum_{i=1}^n \mathbb{1}_{Y_i=1} - \hat{\pi}_1 \sum_{i=1}^n \mathbb{1}_{Y_i=1} = n\hat{\pi}_1 - \hat{\pi}_1 \sum_{i=1}^n \mathbb{1}_{Y_i=1}$$

ipp

$$\boxed{\hat{\pi}_1 = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{Y_i=1}}$$

de plus,

$$\frac{\partial^2 \log \prod_{i=1}^n \hat{\pi}_1^{Y_i} (1-\hat{\pi}_1)^{1-Y_i}}{\partial \hat{\pi}_1^2} = \sum_{i=1}^n \left( -\frac{\mathbb{1}_{Y_i=1}}{\hat{\pi}_1^2} - \frac{1-\mathbb{1}_{Y_i=1}}{(1-\hat{\pi}_1)^2} \right)$$

Ainsi la fonction à optimiser étant concave, le point précédent est bien le maximum.

5) Remarquons que  $V_i, \mathbb{1}_{Y_i=1} \sim \mathcal{B}(\pi_i)$ , et  
Loi de Bernoulli,

de plus que  $\frac{\mathbb{1}_{Y_1=1}}{Y_1}, \dots, \frac{\mathbb{1}_{Y_n=1}}{Y_n}$  sont indépendants.

Alors •  $E\left(\frac{1}{n} \sum_{i=1}^n \frac{\mathbb{1}_{Y_i=1}}{Y_i}\right) = \frac{1}{n} \sum_{i=1}^n E\left(\frac{\mathbb{1}_{Y_i=1}}{Y_i}\right) = \pi_1 \Rightarrow$  sans biais

•  $\text{Var}\left(\frac{1}{n} \sum_{i=1}^n \frac{\mathbb{1}_{Y_i=1}}{Y_i}\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}\left(\frac{\mathbb{1}_{Y_i=1}}{Y_i}\right) = \frac{\pi_1(1-\pi_1)}{n}$



6) Nous allons appliquer l'inégalité de Hoeffding avec

$$V_i, Z_i = \pm 1, \dots, b, 1 \text{ et } a = 0$$

et  $t = \epsilon n$  avec  $\epsilon > 0$

$$\text{Il vient } \mathbb{P}(|\hat{\pi} - \pi| > \epsilon) \leq 2 \exp(-2n\epsilon^2)$$

Ainsi, cet estimateur converge rapidement vers l'estimateur de sa moyenne.

7)  $V_x, p_h(x) = \frac{1}{nb^d} \sum_{i=1}^n \underbrace{\kappa\left(\frac{x-z_i}{h}\right)}_{\geq 0 \text{ par hypothèse}}$

donc  $p_h(\cdot) \geq 0$ .

$$\begin{aligned} \int_{\mathbb{R}^d} p_h(x) dx &= \int_{\mathbb{R}^d} \frac{1}{nb^d} \sum_{i=1}^n \kappa\left(\frac{x-z_i}{h}\right) dx \\ &= \frac{1}{nb^d} \sum_{i=1}^n \int_{\mathbb{R}^d} \kappa\left(\frac{x-z_i}{h}\right) dx \end{aligned}$$

Dans chaque terme de la somme, nous effectuons le changement de variable  $v_i = \frac{x-z_i}{h}$

$$\text{PP } \frac{1}{\pi} (1 - |v_i|)^2 \dots$$

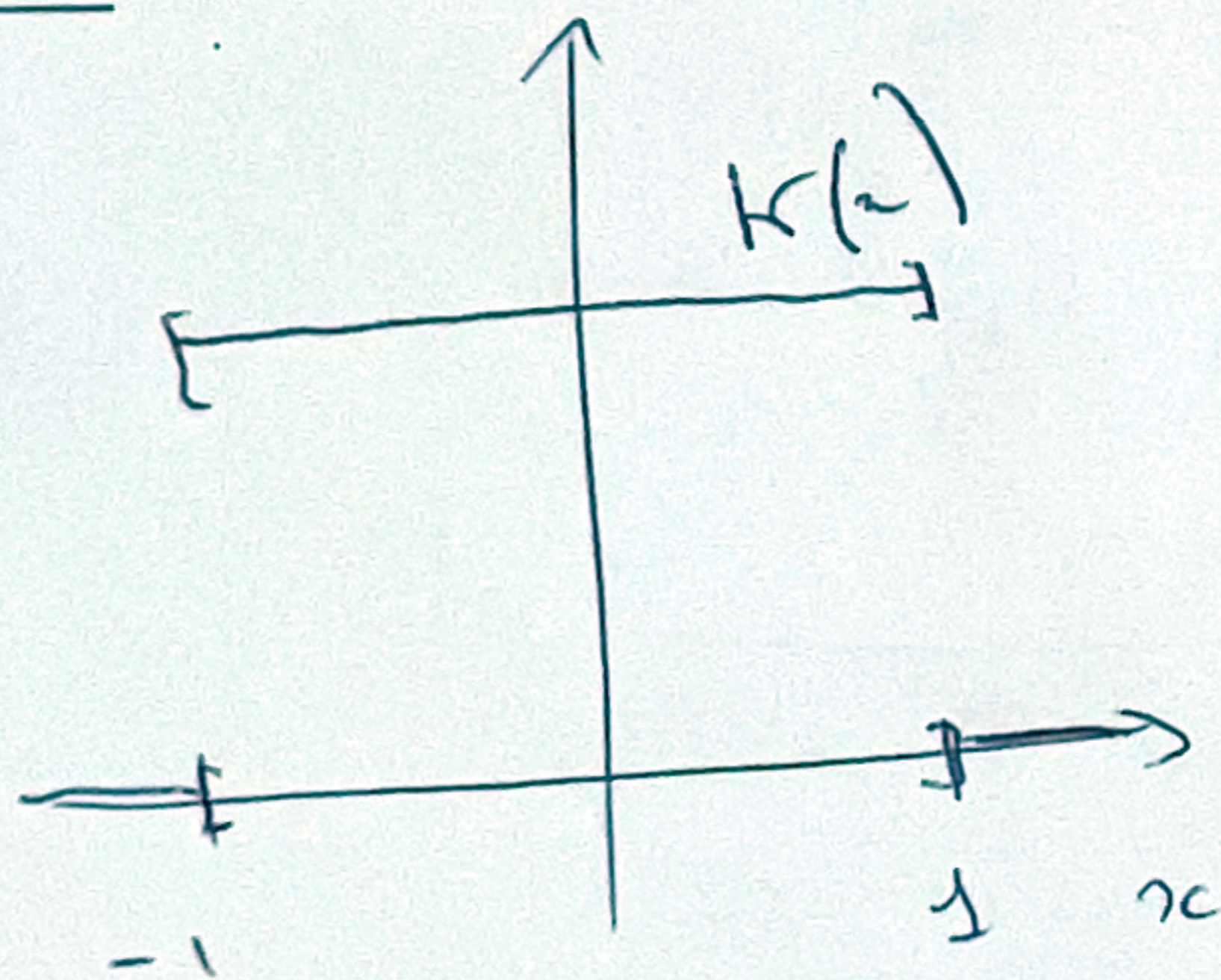


$$\int_{\mathbb{M}^d} p_a(\cdot) du = \frac{1}{n p_d} \sum_{i=1}^n \int_{\mathbb{M}^d} k(u_i) \underbrace{p^d}_{\text{Jacobien du changement de variable}} du_i$$

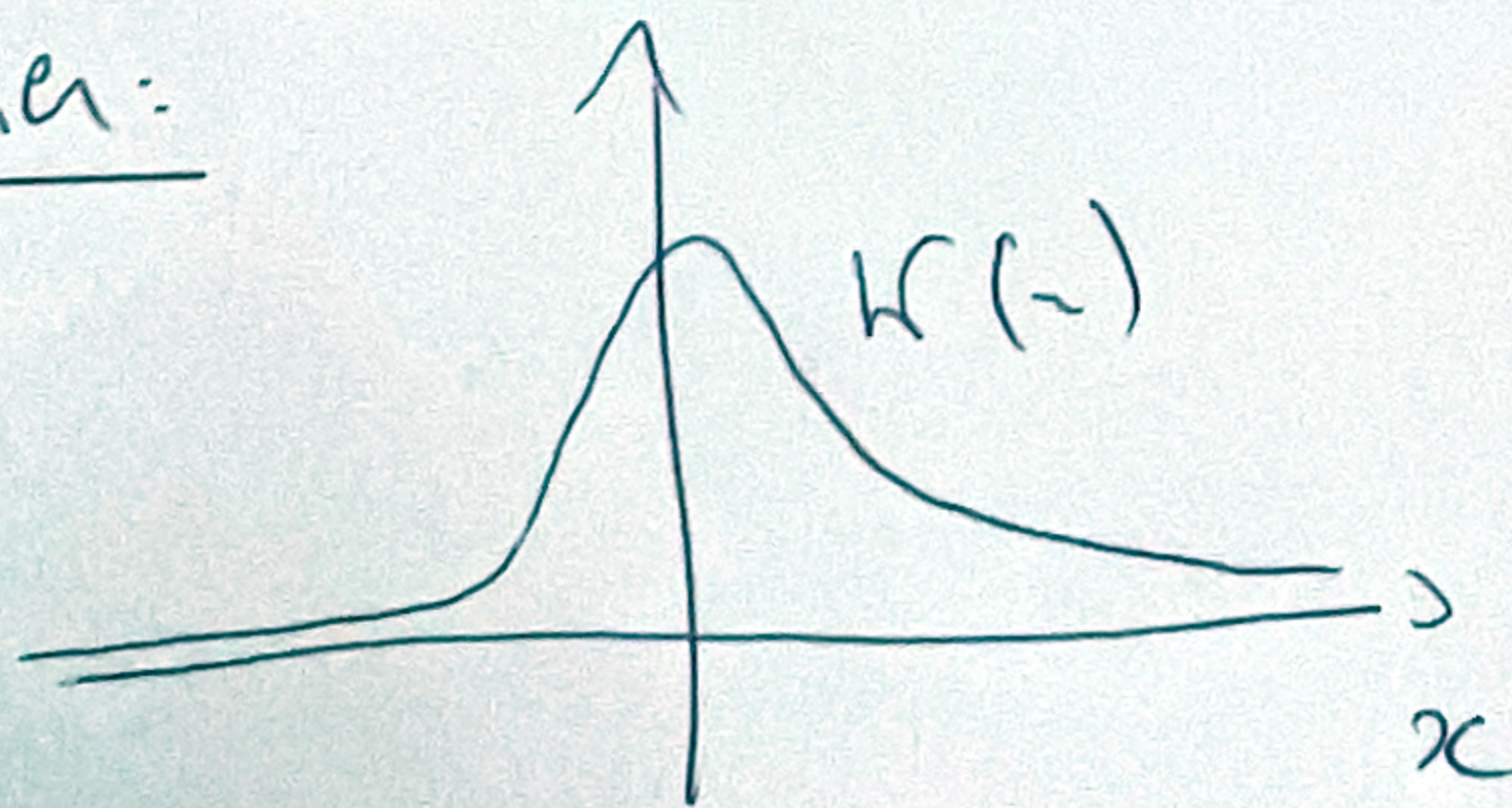
$$= \frac{1}{n} \sum_{i=1}^n \underbrace{\int_{\mathbb{M}^d} k}_{=1}$$

$$= 1.$$

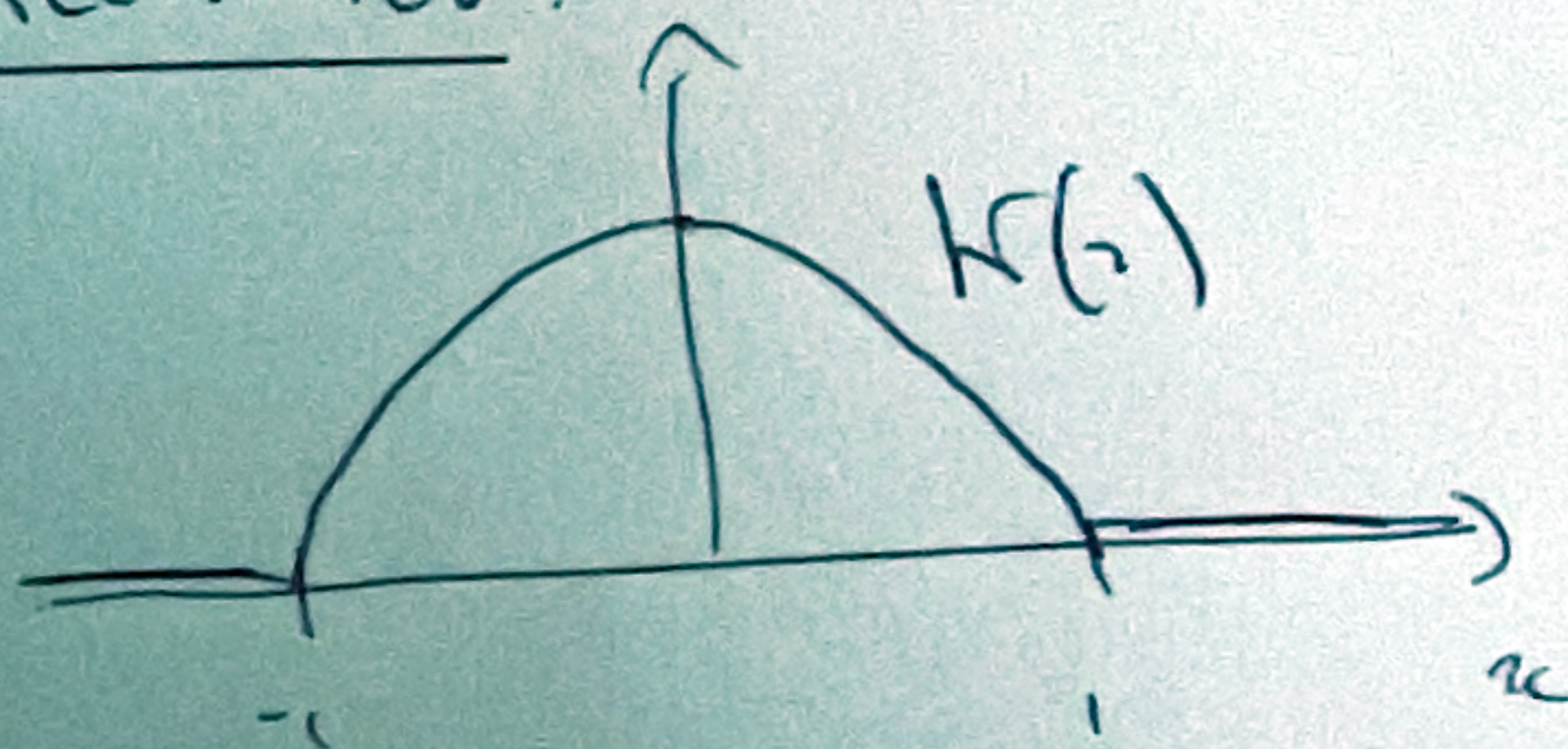
8) Uniforme:



Gaussien:



Épanechnikov:





• Lorsque  $h$  est grand :

Localement, beaucoup d'information est agrégée, il y a peu de variance mais un fort biais.

• Lorsque  $h$  est petit :

Un point n'ajoute de la "masse" que très peu, il y a peu de biais mais une forte variance.

g)

Bon  
Compromis

Sous  
Apprentissage

Sur  
Apprentissage

Le biais est trop grand et l'algorithme n'est pas assez expressif par rapport à la "vraie" fonction de ~~classification~~ classification

La variance est trop importante et l'algorithme est capable d'apprendre le bruit dans les données.



## Exercice 2:

1) Soit  $n \geq 1$ ,  $x_1, \dots, x_n \in X$ ,  $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ ,

$$\begin{aligned} \sum_{i,j} \alpha_i \alpha_j P(x_i, x_j) &= \sum_{i,j} \alpha_i \alpha_j \langle \phi(x_i), \phi(x_j) \rangle \quad (\text{par définition}) \\ &= \left\langle \sum_i \alpha_i \phi(x_i), \sum_j \alpha_j \phi(x_j) \right\rangle \quad (\text{bilinearité}). \\ &= \left\| \sum_i \alpha_i \phi(x_i) \right\|^2 \quad (\text{définition de la norme}) \\ &\geq 0 \quad (\text{car carré d'un positif}). \end{aligned}$$

2) D'après le théorème d'Aronszajn, si  $P$  est un noyau symétrique positif, alors il existe un espace de Hilbert et  $\phi: X \rightarrow H$  tel que

$$\forall x_1, x_2 \in X, P(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle.$$

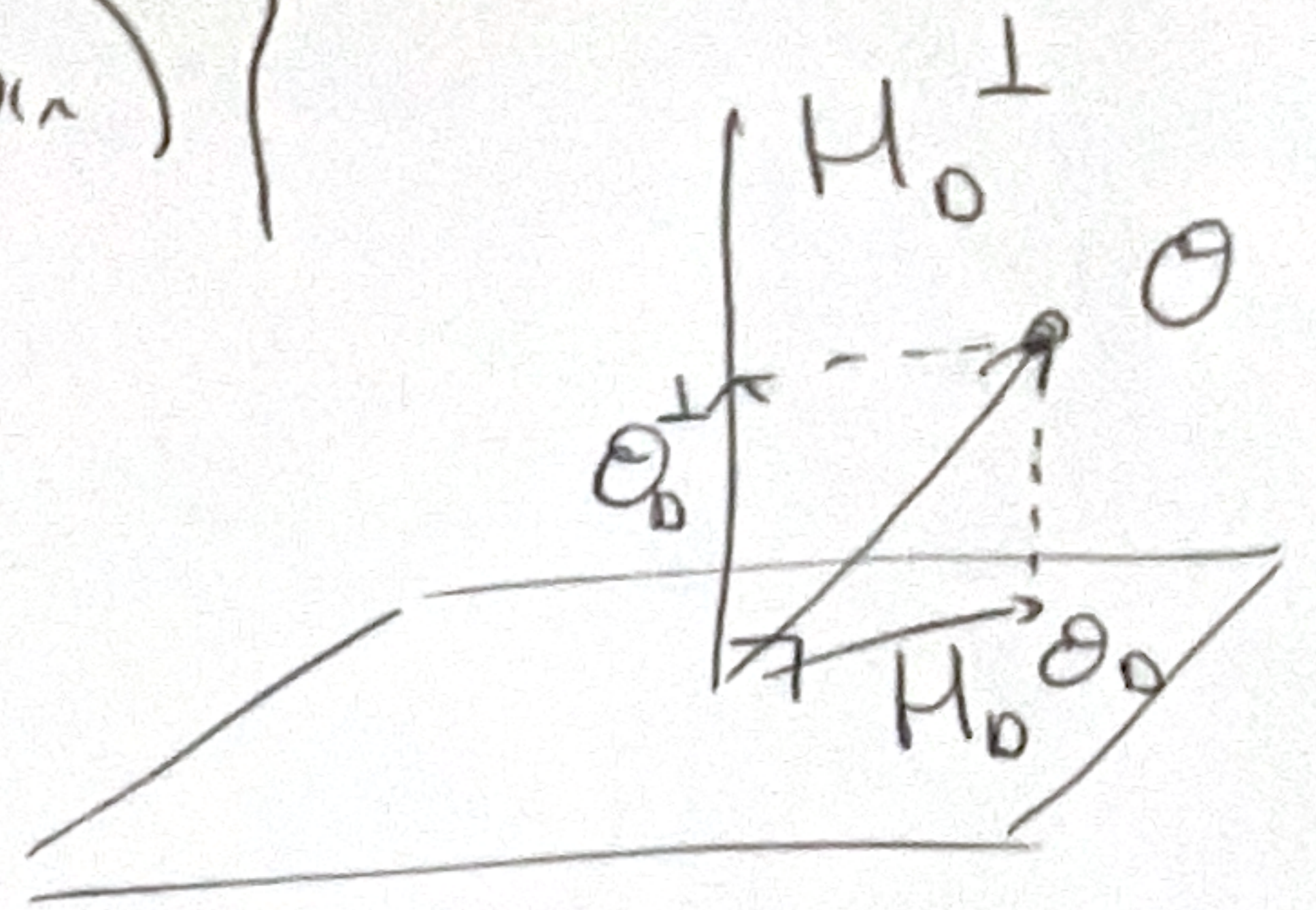
Donc, si on a un noyau pas un produit scalaire et un foncteur  $\phi$  n'est pas restrictif, il y a équivalence entre les deux.



3) Notons  $H_0 \equiv \text{Vect} \{ \phi(x_1), \dots, \phi(x_n) \}$

Sat  $\theta \in H$ ,  $\theta = \theta_0 + \theta_0^\perp$

$\uparrow$                      $\uparrow$   
 $H_0$                  $H_0^\perp$



Notons  $M_n(\theta) \equiv \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-y_i \langle \theta, \phi(x_i) \rangle} \right) + \lambda \|\theta\|^2$

Alors  $M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-y_i \langle \theta_0 + \theta_0^\perp, \phi(x_i) \rangle} \right) + \lambda \|\theta_0 + \theta_0^\perp\|^2$

$= \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-y_i \left( \langle \theta_0, \phi(x_i) \rangle + \langle \theta_0^\perp, \phi(x_i) \rangle \right)} \right) + \lambda \|\theta_0\|^2 + \lambda \|\theta_0^\perp\|^2$

par linéarité

par Pythagore.

or,  $\forall i, \langle \theta_0^\perp, \phi(x_i) \rangle = 0$  car  $\theta_0^\perp \in \text{Vect} \{ \phi(x_1), \dots, \phi(x_n) \}^\perp$ .

donc  $M_n(\theta) = M_n(\theta_0) + \lambda \|\theta_0^\perp\|^2$

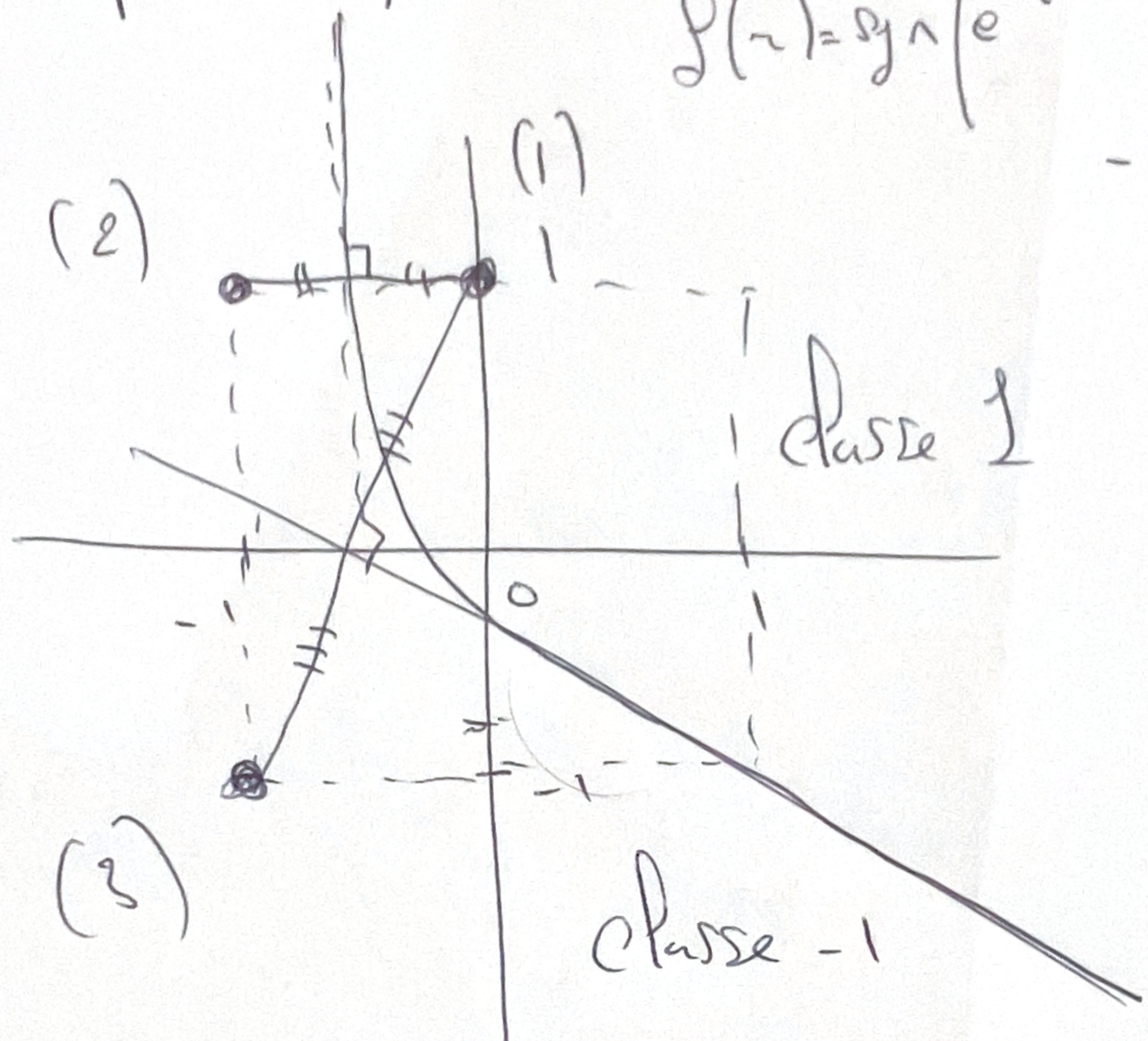
Une solution optimale (si elle existe) a donc une composante dans l'orthogonal de  $H_0$  nulle. On peut se restreindre à chercher dans  $H_0$ .



Sur  $H_3$ , le problème est fortement convexe (grâce au terme  $\lambda \|0\|^2$ ). Il admet donc une unique solution.

1) (noyau exponentiel)

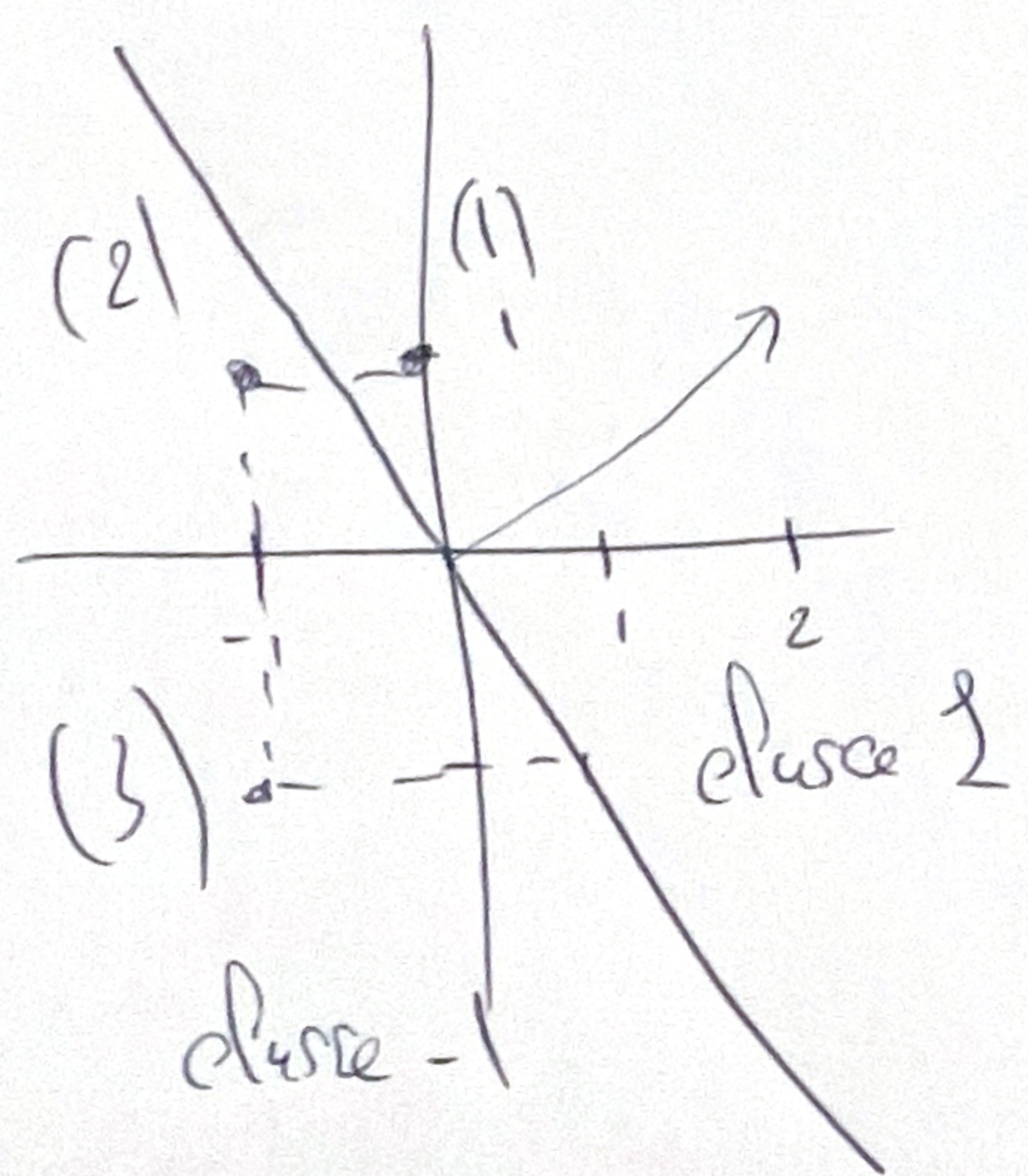
$$f(x) = \text{sgn} \left( e^{-\|x - \begin{pmatrix} 0 \\ 1 \end{pmatrix}\|^2} - \|x - \begin{pmatrix} -1 \\ 1 \end{pmatrix}\|^2 - e^{-\|x - \begin{pmatrix} -1 \\ -1 \end{pmatrix}\|^2} \right)$$



(noyau linéaire)

$$f(x) = \text{sgn} \left( \langle x, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \rangle - \langle x, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \rangle - \langle x, \begin{pmatrix} -1 \\ -1 \end{pmatrix} \rangle \right)$$

$$= \text{sgn} \left( \langle x, \begin{pmatrix} 2 \\ 1 \end{pmatrix} \rangle \right)$$





5) Changement de variables  $\Theta = \alpha, \phi(x_1) + \dots + \alpha_n \phi(x_n)$

$$M_n(\alpha_1, \dots, \alpha_n) = \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-y_i \left( \sum_j \alpha_j \phi(x_j), \phi(x_i) \right)} \right) + \lambda \left\| \sum_j \alpha_j \phi(x_j) \right\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-y_i \sum_j \alpha_j (\phi(x_j), \phi(x_i))} \right) + \lambda \left( \sum_i \alpha_i \phi(x_i), \sum_j \alpha_j \phi(x_j) \right)$$

$$= \frac{1}{n} \sum_{i=1}^n \log \left( 1 + e^{-y_i (K\alpha)_i} \right) + \lambda \alpha^T K \alpha$$

où  $K = \left( \phi(x_i, x_j) \right)_{i,j}$

6) IP est possible de résoudre le problème par descente de gradient (stochastique).

IP peut être reformulé de ce côté

$$\nabla_{\alpha} (\lambda \alpha^T K \alpha) \text{ et } \nabla_{\alpha} \left( \log \left( 1 + e^{-y_i (K\alpha)_i} \right) \right) \forall i$$

ce qui donne

- $\nabla_{\alpha} (\lambda \alpha^T K \alpha) = 2\lambda K \alpha$

- $\nabla_{\alpha} \left( \log \left( 1 + e^{-y_i (K\alpha)_i} \right) \right) = \frac{-e^{-y_i (K\alpha)_i}}{1 + e^{-y_i (K\alpha)_i}} y_i \begin{pmatrix} (\phi(x_1), \phi(x_i)) \\ \vdots \\ (\phi(x_n), \phi(x_i)) \end{pmatrix}$