Examen:

Introduction à l'apprentissage automatique

18 novembre 2023

Aucun document n'est autorisé.

Les questions peuvent être traitées de manière indépendante en admettant les résultats des questions précédentes.

Le barème (sur 20 points, auxquels s'ajoutent 3 points bonus) n'est donné qu'à titre indicatif.

Notations

Dans tout le sujet, on notera :

- 1. 1 le vecteur rempli de 1 (de taille adéquate).
- 2. $\mathbb{1}_A = \begin{cases} 1 & \text{si } A \text{ est vrai} \\ 0 & \text{sinon.} \end{cases}$
- 3. card (I) le cardinal de tout ensemble $I \subset \mathbb{N}$.
- 4. \mathbf{I}_n la matrice identité de taille n (la taille peut varier).
- 5. $\operatorname{sign}: x \in \mathbb{R} \mapsto \begin{cases} 1 & \text{si } x > 0 \\ -1 & \text{sinon.} \end{cases}$
- 6. $\|A\|_F = \sqrt{\operatorname{tr}(AA^\top)}$ la norme de Frobénius de toute matrice A.

Exercice 1 (Algorithme EM, 3½ points)

Soient $(X_1, Y_1), \ldots, (X_n, Y_n)$ des couples i.i.d. à valeurs dans $\mathbb{N} \times \{0, 1\}$ telles que

$$\begin{cases} Y_1 \sim \mathcal{B}(\pi^*), & \pi^* \in]0, 1[\\ X_1 \mid Y_1 \sim \mathcal{P}(\lambda_{Y_1}^*), & \lambda_{Y_1}^* > 0, \end{cases}$$

où $B(\pi^*)$ est la loi de Bernoulli de paramètre π^* et $\mathcal{P}(\lambda)$ la loi de Poisson de paramètre $\lambda > 0$, de densité $x \in \mathbb{N} \mapsto \frac{\lambda^x}{x!} e^{-\lambda}$ par rapport à la mesure de comptage sur \mathbb{N} . Dans la suite, on souhaite partitionner X_1, \ldots, X_n via l'algorithme EM.

- 1. (1 point) On suppose d'abord observer $(X_1, Y_1), \ldots, (X_n, Y_n)$. Définir, pour tout $\theta = (\pi, \lambda_1, \lambda_0) \in]0, 1[\times \mathbb{R}_+^* \times \mathbb{R}_+^*]$ la log-vraisemblance $\ell(\theta, X_1, \ldots, X_n, Y_1, \ldots, Y_n)$ de θ au regard des observations et donner l'estimateur du maximum de vraisemblance de $\theta^* = (\pi^*, \lambda_1^*, \lambda_0^*)$.
- 2. (½ point) Déterminer la loi de $Y_1 \mid X_1$, notée Q_{θ^*,X_1} .

- 3. (1 point) On suppose maintenant n'observer que X_1, \ldots, X_n mais disposer d'un estimateur candidat $\hat{\theta}_0$. Soient alors Z_1, \ldots, Z_n telles que $Z_1 \mid \hat{\theta}_0, \ldots, Z_n \mid \hat{\theta}_0$ sont i.i.d. et pour tout $i \in [1, n]$, $Z_i \mid (X_1, \ldots, X_n) = Z_i \mid (X_i, \hat{\theta}_0) \sim Q_{\hat{\theta}_0, X_i}$. Déterminer $\arg \max_{\theta \in [0, 1[\times \mathbb{R}^*_+ \times \mathbb{R}^*_+]} F(\theta \mid \hat{\theta}_0) = \mathbb{E} \left[\ell(\theta, X_1, \ldots, X_n, Z_1, \ldots, Z_n \mid (X_1, \ldots, X_n) \right]$.
- 4. (1 point) Décrire l'algorithme EM produisant la suite d'estimateurs $(\hat{\theta}_t)_{t\geq 1}$ dans ce cas.
- 5. (1 point (bonus)) En raisonnant de manière générale et en appelant m_{θ^*} la densité marginale de $\mathbf{X} = (X_1, \dots, X_n)$, montrer qu'à chaque iteration $t \geq 1$,

$$\log\left(m_{\hat{\theta}_{t+1}}(\mathbf{X})\right) - \log\left(m_{\hat{\theta}_t}(\mathbf{X})\right) \ge 0.$$

On devra faire intervenir un vecteur aléatoire $\mathbf{Z} = (Z_1, \dots, Z_n)$ tel que

$$\mathbf{Z} \mid \mathbf{X} \sim Q_{\hat{\theta}_t, X_1} \otimes \cdots \otimes Q_{\hat{\theta}_t, X_n}.$$

Exercice 2 (Clustering spectral, $3\frac{1}{2}$ points)

Soient $\{x_1, \ldots, x_n\} \subset \mathbb{R}^d$ un jeu de données, $s : \mathbb{R}^d \times \mathbb{R}^d \to [0, 1]$ une mesure de similarité symétrique et $W = (s(x_i, x_j))_{1 \leq i, j \leq n}$ la matrice de similarité correspondante. Nous avons montré en cours que, pour n'importe quelle partition (I_1, \ldots, I_k) des indices [1, n]

RatioCut
$$(I_1, \ldots, I_k) = \operatorname{tr}(H^{\top}LH),$$

où
$$L = (\left[\sum_{\ell=1}^{n} W_{i,\ell}\right] \mathbb{1}_{i=j} - W_{i,j})_{1 \le i,j \le n}$$
 et $H = \left(\frac{\mathbb{1}_{i \in I_j}}{\sqrt{\operatorname{card}(I_j)}}\right)_{\substack{1 \le i \le n \\ 1 \le j \le k}}$.

- 1. (1 point) Rappeler l'algorithme de clustering spectral construit sur le RatioCut.
- 2. (1 point) Soit $J = (\mathbb{1}_{i \in I_j})_{\substack{1 \leq i \leq n \\ 1 \leq j \leq k}}$. Exprimer la matrice de projection orthogonale sur range(J), notée P_J , qui est l'unique matrice de taille $n \times n$ telle que pour tout $x \in \mathbb{R}^n$,

$${P_J x} = \arg\min_{y \in \operatorname{range}(J)} ||x - y||_{\ell_2}$$

et montrer que sa décomposition en éléments propres est définie par les vecteurs propres $u_j = \left(\frac{\mathbb{1}_{i \in I_j}}{\sqrt{\operatorname{card}(I_j)}}\right)_{1 \leq i \leq n}, \ j = 1 \dots k,$ de valeur propre commune 1.

3. $(1\frac{1}{2} \text{ points})$ Soit $\mathbb{S} = \{U \in \mathbb{R}^{n \times k} : U^{\top}U = \mathbf{I}_k\}$. On souhaite illustrer l'assertion $\ll \text{range}(H)$ est l'espace vectoriel le plus proche de range(J) parmi ceux engendrés par les matrices de l'ensemble $\mathbb{S} \gg$. Pour ce faire, montrer que

$$H \in \operatorname{arg\,min}_{U \in \mathbb{S}} \|P_J - P_U\|_F$$
,

où P_U est le projecteur orthogonale sur range(U).

Exercice 3 (Fonctions de perte pour la classification, 5½ points)

On considère un couple aléatoire (X,Y) à valeurs dans $\mathbb{R}^d \times \{\pm 1\}$ tel que $\eta : x \in \mathbb{R}^d \mapsto \mathbb{P}(Y=1 \mid X=x) \in]0,1[$. Soit maintenant $\varphi : \mathbb{R} \to \mathbb{R}$ une fonction convexe et différentiable autour de 0 telle que

$$\varphi'(0) < 0$$
 et $\arg\min_{u \in \mathbb{R}} \varphi(u) \neq \emptyset$.

Remarquons qu'en particulier :

$$\forall u \in \mathbb{R} : \quad \varphi(u) > \varphi(0) + \varphi'(0)u.$$

1. (1½ points) En remarquant que $\varphi'(0) = \lim_{\substack{u \to 0 \ u > 0}} \frac{\varphi(u) - \varphi(0)}{u}$, montrer que $\exists \bar{u} > 0$: $\varphi(0) > \varphi(\bar{u})$. En déduire que $\forall u < 0, \varphi(u) > \varphi(\bar{u})$, puis que

$$\operatorname{arg\,min}_{u \in \mathbb{R}} \varphi(u) \subset \mathbb{R}_+^*.$$

2. $(1\frac{1}{2} \text{ points})$ Soient $\ell : \mathbb{R} \to \mathbb{R}$ une fonction de perte convexe, positive et différentiable autour de 0 telle que $\ell'(0) < 0$. Justifier que pour tout $x \in \mathbb{R}^d$

$$\varphi: u \in \mathbb{R} \mapsto \mathbb{E}\left[\ell(Yu) \mid X = x\right]$$

est coercive (i.e. $\lim_{-\infty} \varphi = \lim_{\infty} \varphi = \infty$). En appelant u_x^* un miniseur de φ sur \mathbb{R} , montrer que $f^*: x \mapsto u_x^*$ est minimiseur du risque

$$f \mapsto \mathbb{E}[\ell(Yf(X))]$$

et que $g^*: x \mapsto \text{sign}(f^*(x))$ est un classifieur de Bayes.

- 3. (1 point) On choisit $\ell: u \mapsto \max(0, 1-u)^2$. Dessiner le graphe de φ puis exprimer f^* dans ce cas.
- 4. (1½ points) Proposer une manière d'estimer g^* par un classifieur linéaire fondé sur la minimisation d'un risque régularisé construit sur la perte ℓ et exprimer le gradient de ce risque.
- 5. (1 point (bonus)) Montrer que pour $\ell: u \mapsto \frac{1}{(1+e^u)^2}, f^*: x \mapsto \log\left(\frac{\eta(x)}{1-\eta(x)}\right)$. Conclure.

Exercice 4 (Classification à noyau, 7½ points)

Soient $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^d \times \{\pm 1\}$ un jeu de données, $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ un noyau, \mathcal{H} le RKHS associé et $\lambda > 0$. On s'intéresse à la construction d'un classifieur via la résolution du problème d'optimisation

$$\underset{h \in \mathcal{H}, \xi \in \mathbb{R}^n}{\text{minimiser}} \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{2} \sum_{i=1}^n \xi_i^2$$
s. c.
$$\forall i \in [1, n] \begin{cases} Y_i h(X_i) \ge 1 - \xi_i &: \alpha_i \ge 0 \\ \xi_i \ge 0 &: \beta_i \ge 0. \end{cases}$$
(P1)

dans le dual (on donne dans (P1) les multiplicateurs de Lagrange α_i et β_i associés à chaque contrainte).

1. (1 point) Les conditions de qualification de Slater sont-elles vérifiées? Définir le lagrangien \mathcal{L} associé à (P1) et expliciter, pour tous $\alpha \in \mathbb{R}^n_+$ et $\beta \in \mathbb{R}^n_+$, les conditions de stationarité primale en $(h, \xi) \in \mathcal{H} \times \mathbb{R}^n$:

$$\nabla_h \mathcal{L}(h, \xi, \alpha, \beta) = 0$$
 et $\nabla_{\xi} \mathcal{L}(h, \xi, \alpha, \beta) = 0$.

2. $(1\frac{1}{2} \text{ points})$ Montrer qu'un problème dual à (P1) est

$$\underset{\alpha \in \mathbb{R}_{+}^{n}, \, \beta \in \mathbb{R}_{+}^{n}}{\operatorname{maximiser}} \ -\frac{1}{2\lambda} \alpha^{\top} Q \alpha + \mathbb{1}^{\top} \alpha - \frac{1}{2} \left\| \alpha + \beta \right\|_{\ell_{2}}^{2},$$

où Q est une matrice à préciser, puis que ce problème est équivalent (au sens où connaissant les solutions de l'un, on peut déterminer celles de l'autre et vice versa) à

$$\underset{\alpha \in \mathbb{R}^n_+}{\text{minimiser}} \ \frac{1}{2\lambda} \alpha^\top P \alpha - \mathbb{1}^\top \alpha, \tag{P2}$$

où P est une matrice à préciser.

- 3. (1 point) Montrer que P est symétrique et semi-définie positive. Que peut-on en déduire de (P2)?
- 4. (1 point) Expliciter les étapes d'un algorithme de résolution de (P2) de type « descente par coordonnée ».
- 5. $(1\frac{1}{2} \text{ points})$ Énoncer les conditions KKT pour des candidats solutions (h^*, ξ^*) et (α^*, β^*) et en déduire une classifieur issu de la résolution de (P1).
- 6. $(1\frac{1}{2} \text{ points})$ Justifier que, connaissant h^* , on peut choisir $\xi_i^* = \max(0, 1 Y_i h^*(x_i))$ pour tout $i \in [1, n]$. En déduire que pour tout $i \in [1, n]$, si $Y_i h^*(x_i) > 1$, alors $\alpha_i^* = 0$.
- 7. (1 point (bonus)) Proposer un critère d'arrêt pour l'algorithme itératif de la question 4. Le justifier rapidement et expliciter le calcul.

Exercice 1 (Algorithme EM, 3½ points)

Soient $(X_1, Y_1), \ldots, (X_n, Y_n)$ des couples i.i.d. à valeurs dans $\mathbb{N} \times \{0, 1\}$ telles que

$$\begin{cases} Y_1 \sim \mathcal{B}(\pi^*), & \pi^* \in]0, 1[\\ X_1 \mid Y_1 \sim \mathcal{P}(\lambda_{Y_1}^*), & \lambda_{Y_1}^* > 0, \end{cases}$$

où $B(\pi^*)$ est la loi de Bernoulli de paramètre π^* et $\mathcal{P}(\lambda)$ la loi de Poisson de paramètre $\lambda > 0$, de densité $x \in \mathbb{N} \mapsto \frac{\lambda^x}{x!} e^{-\lambda}$ par rapport à la mesure de comptage sur \mathbb{N} . Dans la suite, on souhaite partitionner X_1, \ldots, X_n via l'algorithme EM.

1. (1 point) On suppose d'abord observer $(X_1, Y_1), \ldots, (X_n, Y_n)$. Définir, pour tout $\theta = (\pi, \lambda_1, \lambda_0) \in]0, 1[\times \mathbb{R}_+^* \times \mathbb{R}_+^*]$ la log-vraisemblance $\ell(\theta, X_1, \ldots, X_n, Y_1, \ldots, Y_n)$ de θ au regard des observations et donner l'estimateur du maximum de vraisemblance de $\theta^* = (\pi^*, \lambda_1^*, \lambda_0^*)$.

(i)
$$\forall i \in (1,n)$$
, $P(\gamma_{i=1}) = \lambda$ $P(\gamma_{i=0}) = 1-\lambda$
 $P(\chi_{i=1}) = \frac{\lambda_{i}}{\chi!} e^{-\lambda_{i}} P(\chi_{i=1}) = \frac{\lambda_{o}}{\chi!} e^{-\lambda_{o}}$

donc
$$f_{(x,\gamma)}(x,y) = \left(\frac{\lambda_{x}^{(x)}}{x!} \cdot e^{-\lambda_{1}} \right)^{\left\{ y=1 \right\}} \cdot \left((1-\lambda) \cdot \frac{\lambda_{0}^{(x)}}{x!} \cdot e^{-\lambda_{0}} \right)^{\left\{ y=0 \right\}}$$

$$\log f_{(x,\gamma)}(x,y) = 1_{\left\{ y=1 \right\}} \cdot \left(\log x + x \cdot \log \lambda_{1} - \log x! - \lambda_{1} \right)$$

$$+ 1_{\left\{ y=0 \right\}} \cdot \left(\log (1-\lambda) + x \cdot \log \lambda_{0} - \log x! - \lambda_{0} \right)$$

donc

$$\frac{\left((0,X_{1},...,X_{n},Y_{1},...,Y_{n}) = \sum_{i=1}^{n} \mathbb{1}_{\{Y_{i}=1\}} \cdot (\log x + X_{i} \cdot \log x_{1} - \log X_{i}! - \lambda_{1})}{+ \sum_{i=1}^{n} \mathbb{1}_{\{Y_{i}=0\}} \cdot (\log (1-x) + X_{i} \cdot \log x_{0} - \log X_{i}! - \lambda_{0})}$$

$$\frac{\partial \ell}{\partial \lambda} = 0 \implies \sum_{i=1}^{n} 1_{\{i_{i}=1\}} \cdot \frac{1}{\lambda} - \sum_{i=1}^{n} 1_{\{i_{i}=0\}} \frac{1}{1-\lambda} = 0$$

$$\lambda^{*} = \frac{1}{n} \sum_{i=1}^{n} 1_{\{i_{i}=1\}}$$

$$\frac{\partial \ell}{\partial \lambda_{i}} = 0 \qquad \Longrightarrow \qquad \sum_{\underline{i}=1}^{n} \mathbb{1}_{\left\{Y_{\underline{i}}=1\right\}} \cdot \left(\frac{X_{\underline{i}}}{\lambda_{i}} - 1\right) = 0$$

$$\lambda_{i}^{+} = \frac{\sum_{\underline{i}=1}^{n} \mathbb{1}_{\left\{Y_{\underline{i}}=1\right\}} \cdot X_{\underline{i}}}{\sum_{\underline{i}=1}^{n} \mathbb{1}_{\left\{Y_{\underline{i}}=1\right\}}}$$

$$\frac{\partial \ell}{\partial \lambda_o} = 0 \qquad \Longrightarrow \qquad \lambda_o^* = \frac{\sum_{i=1}^n \mathbb{1}_{\{Y_i = 0\}} \cdot X_i}{\sum_{i=1}^n \mathbb{1}_{\{Y_i = 0\}}}$$

2. (½ point) Déterminer la loi de $Y_1 \mid X_1$, notée Q_{θ^*,X_1} .

$$P(Y_{i=1}|X_{i}) = \frac{P(Y_{i=1},X_{i}=x)}{P(X_{i}=x)} = \frac{z \cdot \lambda_{i}^{X_{i}} e^{-\lambda_{i}}}{z \cdot \lambda_{i}^{X_{i}} e^{-\lambda_{i}} + (1-z) \cdot \lambda_{o}^{X_{i}} e^{-\lambda_{o}}} = q_{0,X_{i}}$$

$$Y_{i}|X_{i} \sim B(\frac{z_{i} \lambda_{i}^{x_{i}} e^{-\lambda_{i}}}{z_{i} \lambda_{i}^{x_{i}} e^{-\lambda_{i}} + (1-z_{i}) \cdot \lambda_{o}^{x_{i}} e^{-\lambda_{o}}})$$

3. (1 point) On suppose maintenant n'observer que X_1, \ldots, X_n mais disposer d'un estimateur candidat $\hat{\theta}_0$. Soient alors Z_1, \ldots, Z_n telles que $Z_1 \mid \hat{\theta}_0, \ldots, Z_n \mid \hat{\theta}_0$ sont i.i.d. et pour tout $i \in [1, n], Z_i \mid (X_1, \ldots, X_n) = Z_i \mid (X_i, \hat{\theta}_0) \sim Q_{\hat{\theta}_0, X_i}$. Déterminer

 $\arg\max_{\theta\in]0,1[\times\mathbb{R}_{+}^{*}\times\mathbb{R}_{+}^{*}} F(\theta \mid \hat{\theta}_{0}) = \mathbb{E}\left[\ell(\theta,X_{1},\ldots,X_{n},Z_{1},\ldots,Z_{n} \mid (X_{1},\ldots,X_{n})\right].$

$$\begin{array}{ll}
|3| & F(\theta|\hat{\theta}_{0}) = \mathbb{E}\left[\sum_{i=1}^{n} 1_{\{Z_{i}=1\}} \cdot (\log x + X_{i} \cdot \log x_{i} - \log X_{i} - \lambda_{i}) \\
& + \sum_{i=1}^{n} 1_{\{Z_{i}=0\}} \cdot (\log (-x) + X_{i} \cdot \log x_{0} - \log X_{i} - \lambda_{0}) | (X_{1},...,X_{n}) \right] \\
& = \sum_{i=1}^{n} 9_{\theta_{0}} \cdot X_{i} \cdot (\log x + X_{i} \cdot \log x_{0} - \log X_{i} - \lambda_{0}) \\
& + \sum_{i=1}^{n} (1 - 9_{\theta_{0}} \cdot X_{i}) \cdot (\log (-x) + X_{i} \cdot \log x_{0} - \log x_{i} - \lambda_{0})
\end{array}$$

donc

$$\frac{\partial F}{\partial \lambda} = 0 \implies \lambda^* = \frac{1}{n} \sum_{i=1}^n \frac{\hat{\lambda}_i \cdot \hat{\lambda}_i^{X_i} \cdot e^{-\hat{\lambda}_i}}{\hat{\lambda}_i \cdot \hat{\lambda}_i^{X_i} \cdot e^{-\hat{\lambda}_i} \cdot \hat{\lambda}_i^{X_i} \cdot e^{-\hat{\lambda}_i}}$$

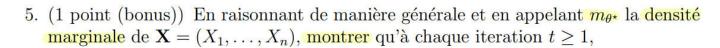
$$\frac{\partial F}{\partial \lambda_{i}} = 0 \implies \lambda_{i}^{*} = \frac{\sum_{i=1}^{n} q_{\hat{\theta}_{0}, X_{i}} \cdot X_{i}}{\sum_{i=1}^{n} q_{\hat{\theta}_{0}, X_{i}}}$$

$$\frac{\partial F}{\partial \lambda_{o}} = 0 \implies \lambda_{o}^{*} = \frac{\sum_{i=1}^{n} (1 - q_{\widehat{\theta}_{o}, X_{i}}) \cdot \chi_{i}}{\sum_{i=1}^{n} (1 - q_{\widehat{\theta}_{o}, X_{i}})}$$

(4) <u>L'algorithme</u> EM
Input: X1,, Xn
Initialiser $\hat{\theta}_o = (\hat{\lambda}_o, \hat{\lambda}_o^{(o)}, \hat{\lambda}_o^{(o)})$
pour le i'ème étape: $\hat{\lambda}_{i} \leftarrow \frac{1}{n} \sum_{i=1}^{n} q_{\hat{\theta}_{i-1}, X_{i}}$
$\hat{\lambda}_i \leftarrow \frac{1}{n} \sum_{i=1}^n q_{\hat{\theta}_{i-1}, X_i}$
$ \hat{\sum}_{i=1}^{n} q_{\hat{\theta}_{i-1}, X_{i}} \times \sum_{i=1}^{n} q_{\hat{\theta}_{i-1}, X_{i}} $
$\sum_{i=1}^{n} q_{\hat{\theta}_{i-1}, X_{i}}$
$ \hat{\lambda}_{0}^{(i)} \leftarrow \frac{\sum_{i=1}^{N} (1 - 9\hat{\theta}_{i-1}, X_{i}) \cdot X_{i}}{\sum_{i=1}^{N} (1 - 9\hat{\theta}_{i-1}, X_{i})} $
$\sum_{i=1}^{k} \left(1 - 9\hat{\theta}_{i-1}, X_{i}\right)$
Output 9êz,x.,, 9êz,xn

4. (1 point) Décrire l'algorithme EM produisant la suite d'estimateurs $(\hat{\theta}_t)_{t\geq 1}$ dans ce

cas.



$$\log \left(m_{\hat{\theta}_{t+1}}(\mathbf{X}) \right) - \log \left(m_{\hat{\theta}_t}(\mathbf{X}) \right) \ge 0.$$

On devra faire intervenir un vecteur aléatoire $\mathbf{Z} = (Z_1, \dots, Z_n)$ tel que

$$\mathbf{Z} \mid \mathbf{X} \sim Q_{\hat{\theta}_t, X_1} \otimes \cdots \otimes Q_{\hat{\theta}_t, X_n}.$$

(5)
$$P(X=x) = \lambda \cdot \frac{\lambda_1^{x}}{x!} \cdot e^{-\lambda_1} + (1-\lambda_1) \cdot \frac{\lambda_0^{x}}{x!} \cdot e^{-\lambda_0}$$

Donc $m_{\widehat{\theta}_{k}}(X) = \prod_{i=1}^{n} \left(q_{\widehat{\theta}_{k}, X_{i}} \cdot \frac{\lambda_{i}^{X_{i}}}{X_{i}!} \cdot e^{\lambda_{i}} + \left(1 - q_{\widehat{\theta}_{k}, X_{i}} \right) \cdot \frac{\lambda_{o}^{X_{i}}}{X_{i}!} \cdot e^{\lambda_{o}} \right)$ $\log(m_{\widehat{\theta}_{k}}(X)) = \sum_{i=1}^{n} \left(\log q_{\widehat{\theta}_{k}, X_{i}} + X_{i} \cdot \log \hat{\lambda}_{i}^{(k)} - \log(X_{i}!) - \hat{\lambda}_{i}^{(k)} \right)$ $+ \log\left(1 - q_{\widehat{\theta}_{k}, X_{i}} \right) + X_{i} \cdot \log \hat{\lambda}_{o}^{(k)} - \log(X_{i}!) - \hat{\lambda}_{o}^{(k)} \right)$

$$\log (m_{\theta_{\ell},i}(X)) - \log (m_{\theta_{\ell}}(X)) = \sum_{i=1}^{n} \left(\log \frac{9 \hat{\theta}_{\ell}, X_{i}}{9 \hat{\theta}_{\ell}, X_{i}} + X_{i} \cdot \log \frac{\hat{\lambda}_{i}^{(\ell+1)}}{\hat{\lambda}_{i}^{(\ell+1)}} - (\hat{\lambda}_{i}^{(\ell+1)} - \hat{\lambda}_{i}^{(\ell)}) \right) \\ + \log \frac{1 - 9 \hat{\theta}_{\ell}, X_{i}}{1 - 9 \hat{\theta}_{\ell}, X_{i}} + X_{i} \cdot \log \frac{\hat{\lambda}_{i}^{(\ell+1)}}{\hat{\lambda}_{i}^{(\ell+1)}} - (\hat{\lambda}_{i}^{(\ell+1)} - \hat{\lambda}_{i}^{(\ell)}) \right)$$

- - -

Exercice 2 (Clustering spectral, 3½ points)

Soient $\{x_1, \ldots, x_n\} \subset \mathbb{R}^d$ un jeu de données, $s : \mathbb{R}^d \times \mathbb{R}^d \to [0, 1]$ une mesure de similarité symétrique et $W = (s(x_i, x_j))_{1 \le i, j \le n}$ la matrice de similarité correspondante. Nous avons montré en cours que, pour n'importe quelle partition (I_1, \ldots, I_k) des indices [1, n]

RatioCut
$$(I_1, \ldots, I_k) = \operatorname{tr}(H^{\top}LH),$$

où
$$L = (\left[\sum_{\ell=1}^{n} W_{i,\ell}\right] \mathbb{1}_{i=j} - W_{i,j})_{1 \le i,j \le n}$$
 et $H = \left(\frac{\mathbb{1}_{i \in I_j}}{\sqrt{\operatorname{card}(I_j)}}\right)_{\substack{1 \le i \le n \\ 1 \le j \le k}}$.

1. (1 point) Rappeler l'algorithme de clustering spectral construit sur le RatioCut.

1) L'algorithme de clustering spectral (Ratio C	(ut)
---	------

Unnormalized	Spectral	Clustering
	1	0

Input: $W \in \mathbb{R}^{n \times n}$ (adjency matrix)

2 - Laplacian de W

H = k minor eigenvectors of L as columns

Yi the row of H

 $(\hat{c}_1,...,\hat{c}_k) \leftarrow$ output of k-means algorithm based on $(\gamma_1,...,\gamma_n)$

Output: $(\hat{c}_1,...,\hat{c}_k)$

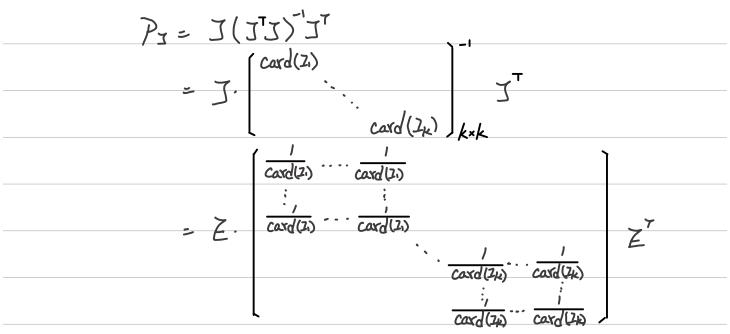
2.	(1 point)	Soit J	= ($\mathbb{1}_{i \in I_j}$	$\leq i \leq n$. Expr	imer la n	natr	rice de	projec	ction of	ortho	gonal	e sur
	range(J),													
	$x \in \mathbb{R}^n$.													

$$\{P_J x\} = \arg\min_{y \in \operatorname{range}(J)} \|x - y\|_{\ell_2}$$

et montrer que sa décomposition en éléments propres est définie par les vecteurs propres $u_j = \left(\frac{\mathbb{1}_{i \in I_j}}{\sqrt{\operatorname{card}(I_j)}}\right)_{1 \le i \le n}, \ j = 1 \dots k,$ de valeur propre commune 1.

(2)		11,ez.	111012	1 _{161k}	÷12				
	J=	:	•	:		3 = E.	;		
		[1 _{ne],}	1 _{nel} ,	1 nelk	n×k			/ : /	7l×,

on a



Done Djelika

3.	(1½ points) Soit $S = \{U \in \mathbb{R}^{n \times k} : U^{\top}U = \mathbf{I}_k\}$. On souhaite illustrer l'assertion
	\ll range (H) est l'espace vectoriel le plus proche de range (J) parmi ceux engendrés
	par les matrices de l'ensemble S ». Pour ce faire, montrer que

$$H \in \operatorname{arg\,min}_{U \in \mathbb{S}} \|P_J - P_U\|_F$$
,

où P_U est le projecteur orthogonale sur $\operatorname{range}(U)$.

argmin $ P_3 - P_0 _F = \operatorname{argmin}_{V \in \mathbb{R}} \sqrt{tr((P_3 - P_0)^T)}$ $V \in S$ $V \in S$
UU=Ix
alors on a $ P_3 - P_0 _{\mathcal{F}} > 0$
comme P4 = P5, alors on a
$\ P_3 - P_A\ _{\mathcal{F}} = 0$
donc $H \in \operatorname{argmin} \ P_3 - P_0 \ _{Z}$

Exercice 3 (Fonctions de perte pour la classification, 5½ points)

On considère un couple aléatoire (X,Y) à valeurs dans $\mathbb{R}^d \times \{\pm 1\}$ tel que $\eta: x \in \mathbb{R}^d \mapsto \mathbb{P}(Y=1 \mid X=x) \in]0,1[$. Soit maintenant $\varphi: \mathbb{R} \to \mathbb{R}$ une fonction convexe et différentiable autour de 0 telle que

$$\varphi'(0) < 0$$
 et $\arg\min_{u \in \mathbb{R}} \varphi(u) \neq \emptyset$.

Remarquons qu'en particulier :

$$\forall u \in \mathbb{R} : \quad \varphi(u) \ge \varphi(0) + \varphi'(0)u.$$

1. (1½ points) En remarquant que $\varphi'(0) = \lim_{\substack{u \to 0 \\ u > 0}} \frac{\varphi(u) - \varphi(0)}{u}$, montrer que $\exists \bar{u} > 0$: $\varphi(0) > \varphi(\bar{u})$. En déduire que $\forall u \leq 0, \varphi(u) > \varphi(\bar{u})$, puis que

$$\operatorname{arg\,min}_{u \in \mathbb{R}} \varphi(u) \subset \mathbb{R}_+^*.$$

11) Par l'hypothèse, aramin (g(u) + \$\phi\$

donc il existe u∈R tel que (g'(u) = 0

et comme 9 est convexe et 9(0) < 0

abors il existe $\overline{u} > 0$, tel que $9'(\overline{u}) = 0$

ed 460> 4(<u>u</u>)

Et comme 900 < 0, 9 est convexe

alors \ u ≤ 0 , \ g(u) < 0

donc 9 est strictement décroissant dans 1-00,0]

donc \u < 0, g(u) > g(v) > g(v)

Donc arginin G(u) C/R+

2. $(1\frac{1}{2} \text{ points})$ Soient $\ell : \mathbb{R} \to \mathbb{R}$ une fonction de perte convexe, positive et différentiable autour de 0 telle que $\ell'(0) < 0$. Justifier que pour tout $x \in \mathbb{R}^d$

$$\varphi: u \in \mathbb{R} \mapsto \mathbb{E}\left[\ell(Yu) \mid X = x\right]$$

est coercive (i.e. $\lim_{-\infty} \varphi = \lim_{\infty} \varphi = \infty$). En appelant u_x^* un miniseur de φ sur \mathbb{R} , montrer que $f^*: x \mapsto u_x^*$ est minimiseur du risque

$$f \mapsto \mathbb{E}[\ell(Yf(X))]$$

et que $g^*: x \mapsto \text{sign}(f^*(x))$ est un classifieur de Bayes.

lim
$$y(u) = \lim_{x \to \infty} \mathbb{E} \left[l(yu) | X = x \right]$$

$$= \mathbb{E} \left[\lim_{x \to +\infty} l(yu) | X = x \right]$$

comme l est convexe, positive et différentiable

alors $\lim_{x \to +\infty} l(u) = \lim_{x \to -\infty} l(u) = +\infty$

donc $\lim_{x \to +\infty} l(yu) = \lim_{x \to -\infty} l(yu) = +\infty$, $\forall y \in \{\pm 1\}$

donc $\lim_{x \to +\infty} l(yu) = \lim_{x \to -\infty} l(yu) = +\infty$
 $\lim_{x \to +\infty} l(x) = \lim_{x \to -\infty} l(y) = +\infty$
 $\lim_{x \to +\infty} l(x) = \lim_{x \to -\infty} l(y) = +\infty$

Comme ux est un miniseur de 4, alors

Comme
$$u_{x} \in \mathcal{E}(\mathcal{X}) = \mathcal{E}[l(Y + f(x))]$$

$$= \mathcal{E}[l(Y + f(x)) | X = x]$$

$$= \mathcal{E}[l(Y \cdot u_{x}^{*}) | X = x]$$

donc f* est un minimiseur du risque

$$S_i \quad \forall = 1$$
, $y(u) = \mathbb{E}[l(u)|X=x]$

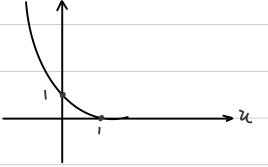
d'après (1)
$$u_x^* \in \mathbb{R}_+$$
 alors $g(x) = 1$

$$S_i = 1$$
, $S(u) = \mathbb{E}[l(-u)|X=x]$

d'après (1)
$$u_x^* \in \mathbb{R}$$
 alors $g(x) = -1$

3. (1 point) On choisit $\ell: u \mapsto \max(0, 1-u)^2$. Dessiner le graphe de φ puis exprimer f^* dans ce cas.

(3) Le graphe



$$R(f) = \mathbb{E}\left[l(Yf(x))\right] = \mathbb{E}\left[max(0, 1-Yf(x))^{2}\right]$$

$$c-a-d$$
 si $Y=1$, $f(x) > 1$

封

7 f(x)=wx+b

4. (1½ points) Proposer une manière d'estimer g^* par un classifieur linéaire fondé sur la minimisation d'un risque régularisé construit sur la perte ℓ et exprimer le gradient de ce risque.

(4) Pour estimer $g^*(x) = sign(f^*(x)) = sign(w^{*T}x + b^*)$

on doit estimer (w*, b*)

On prend N échantillons (X1, Y1), ..., (Xn, Yn)

qui sont i.i.d. de (X,Y)

D les paramétres (w*, b*) est solution de

minimise $\frac{\lambda}{2} \|w\|_{L_x}^2 + \frac{1}{n} \sum_{i=1}^n \max(0, 1-\gamma_i f(x_i))^2$ were

 $\iff \text{minimise } \frac{1}{7} \|w\|_{C_{2}}^{2} + \frac{1}{7n} \sum_{i=1}^{n} \frac{x^{2}}{3i}$ $\text{We} \mathbb{R}^{d}$ $\text{Se} \mathbb{R}$

S.C. $\forall i \in [1,n]$, $\S_i > 0$ $\forall i \in [1,n]$, $\forall i \neq (x_i) > 1-\S_i$

D'après 13), pour f^* , on a $\forall i f^*(x_i) \ge 1$, $\forall i \in [1,n]$

c-d-d == 0, Hielin]

Donc on a le nouveau problème

minimise $\frac{1}{2} \|w\|_{\ell_2}^2$ ber

S.C. Vielin], /¿(w/Xi+b) > 1

Alors on a $\hat{g}(x) = sign(\hat{w}^Tx + \hat{b})$

ona

$$\mathcal{R}_{n}(w,b) = \frac{1}{n} \sum_{i=1}^{n} \max(0, 1 - Y_{i}(w^{T}X_{i} + b))^{2}$$

$$= \frac{1}{n} \sum_{i=1}^{n} (1 - Y_{i}(w^{T}X_{i} + b))^{2} \cdot 1_{\{Y_{i}(w^{T}X_{i} + b) \leq 1\}}$$

$$\nabla_{w} R_{n}(w,b) = -\frac{7}{n} \sum_{i=1}^{n} \langle_{i} X_{i} \cdot (1 - \langle_{i} (w^{T} X_{i} + b)) \cdot 1 | \{ \langle_{i} (w^{T} X_{i} + b) \leq 1 \}$$

$$\nabla_{b} R_{n}(w,b) = -\frac{7}{n} \sum_{i=1}^{n} \langle_{i} \cdot (1 - \langle_{i} (w^{T} X_{i} + b)) \cdot 1 | \{ \langle_{i} (w^{T} X_{i} + b) \leq 1 \}$$

5. (1 point (bonus)) Montrer	que	pour	ℓ	:	u	\mapsto	$\frac{1}{(1+e^u)^2},$	f^{\star}	•	x	\mapsto	log	$\left(\frac{\eta(x)}{1-\eta(x)}\right)$	•
---------------------	-----------	-----	------	--------	---	---	-----------	------------------------	-------------	---	---	-----------	-----	--	---

(5) Pour
$$l(u) = \frac{1}{(1+e^u)^2}$$
 $l'(u) = -\frac{2\cdot e^u}{(1+e^u)^3}$

on a

$$R(f) = \mathbb{E}\left[l(Yf(x))\right]$$

$$= \mathbb{E}\left[\eta(x)\cdot l(f(x)) + (1-\eta(x))\cdot l(-f(x))\right]$$

pour fonction

$$C_{\eta}(u) = \eta \cdot \ell(u) + (1-\eta) \cdot \ell(-u)$$

on a

$$C_{\eta}'(u) = \eta \cdot \hat{l}(u) + (1-\eta) \cdot \hat{l}(-u)$$

$$= -\eta \cdot \frac{2 \cdot e^{u}}{(1+e^{u})^{3}} + (1-\eta) \cdot \frac{2 \cdot e^{u}}{(1+e^{-u})^{3}}$$

alors

$$C_{\eta}(u) = 0 \implies \eta \cdot \frac{2 \cdot e^{u}}{(1 + e^{u})^{3}} = (1 - \eta) \cdot \frac{2 \cdot e^{u}}{(1 + e^{-u})^{3}}$$

$$\frac{\eta \cdot e^{u}}{(1 + e^{u})^{3}} = \frac{(1 - \eta) \cdot e^{2u}}{(1 + e^{u})^{3}}$$

$$e^{u} = \frac{\eta}{1 - \eta}$$

$$u^{*} = \log(\frac{\eta}{1 - \eta})$$

donc

$$f'(x) = \log(\frac{\eta(x)}{1 - \eta(x)})$$

Exercice 4 (Classification à noyau, 7½ points)

Soient $\{(x_1, y_1), \dots, (x_n, y_n)\} \subset \mathbb{R}^d \times \{\pm 1\}$ un jeu de données, $k : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ un noyau, \mathcal{H} le RKHS associé et $\lambda > 0$. On s'intéresse à la construction d'un classifieur via la résolution du problème d'optimisation

$$\underset{h \in \mathcal{H}, \, \xi \in \mathbb{R}^n}{\text{minimiser}} \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 + \frac{1}{2} \sum_{i=1}^n \xi_i^2$$
s. c.
$$\forall i \in [1, n] \begin{cases} Y_i h(X_i) \ge 1 - \xi_i &: \alpha_i \ge 0 \\ \xi_i \ge 0 &: \beta_i \ge 0. \end{cases}$$
(P1)

dans le dual (on donne dans (P1) les multiplicateurs de Lagrange α_i et β_i associés à chaque contrainte).

1. (1 point) Les conditions de qualification de Slater sont-elles vérifiées? Définir le lagrangien \mathcal{L} associé à (P1) et expliciter, pour tous $\alpha \in \mathbb{R}^n_+$ et $\beta \in \mathbb{R}^n_+$, les conditions de stationarité primale en $(h, \xi) \in \mathcal{H} \times \mathbb{R}^n$:

$$\nabla_h \mathcal{L}(h, \xi, \alpha, \beta) = 0$$
 et $\nabla_{\xi} \mathcal{L}(h, \xi, \alpha, \beta) = 0$.

alors
$$\exists \tilde{\xi} = 21$$
, $\hat{h}(x) = 0$, tel que

$$m_i(\tilde{z}) < 0$$

donc les conditions de qualification de Slater sont vérifiés.

Lagrangien L

$$\mathcal{L}(\Lambda, \mathcal{Z}, \alpha, \beta) = \frac{\lambda}{2} \|\Lambda\|_{\mathcal{A}}^{2} + \frac{1}{2} \sum_{i=1}^{n} \mathcal{Z}_{i}^{2} + \sum_{i=1}^{n} \alpha_{i} \cdot (1 - \mathcal{Z}_{i} - \gamma_{i} \cdot \Lambda(x_{i})) - \sum_{i=1}^{n} \beta_{i} \mathcal{Z}_{i}$$

$$\nabla_{\lambda} \mathcal{L}(\lambda, 3, \alpha, \beta) = \lambda \cdot \lambda - \sum_{i=1}^{n} \alpha_{i} \cdot Y_{i} \cdot k(\cdot, X_{i})$$

$$\nabla_{h} \mathcal{L}(h, \mathbf{x}, \alpha, \beta) = 0 \implies \hat{\lambda} = \frac{1}{3} \sum_{i=1}^{n} \alpha_{i} \cdot \gamma_{i} \cdot k(\cdot, X_{i})$$

$$\nabla_{\xi} L(h, \xi, \alpha, \beta) = 0 \implies \hat{\xi} = \alpha + \beta$$

2. (1½ points) Montrer qu'un problème dual à (P1) est

$$\underset{\alpha \in \mathbb{R}_+^n, \, \beta \in \mathbb{R}_+^n}{\operatorname{maximiser}} \ -\frac{1}{2\lambda} \alpha^\top Q \alpha + \mathbb{1}^\top \alpha - \frac{1}{2} \left\| \alpha + \beta \right\|_{\ell_2}^2,$$

où Q est une matrice à préciser, puis que ce problème est équivalent (au sens où connaissant les solutions de l'un, on peut déterminer celles de l'autre et vice versa) à

$$\underset{\alpha \in \mathbb{R}^n_+}{\text{minimiser}} \ \frac{1}{2\lambda} \alpha^\top P \alpha - \mathbb{1}^\top \alpha, \tag{P2}$$

où P est une matrice à préciser.

(2) Par définition

$$D(\alpha,\beta) = \inf_{h \in \mathcal{H}} \mathcal{L}(h, \mathbf{3}, \alpha, \beta) = \mathcal{L}(\hat{h}, \hat{\mathbf{3}}, \alpha, \beta)$$

$$= \frac{\lambda}{2} \cdot \left\| \frac{1}{\lambda} \sum_{i=1}^{n} \alpha_i \cdot \gamma_i \cdot k(\cdot, x_i) \right\|_{\mathcal{H}}^2 + \frac{1}{\lambda} (\alpha + \beta)^{\mathsf{T}} (\alpha + \beta)$$

$$+\sum_{i=1}^{n}\alpha_{i}\cdot\left(1-\alpha_{i}-\beta_{i}-\gamma_{i}\cdot\frac{1}{\beta}\sum_{j=1}^{n}\alpha_{j}\cdot\gamma_{j}\cdot k(x_{i},x_{j})\right)-\beta^{7}(\alpha+\beta)$$

$$= \frac{1}{2\lambda} \cdot \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{j} \lambda_{i} \lambda_{j} \cdot k(X_{i}, X_{j}) - \frac{1}{2} \|\alpha + \beta\|_{\ell_{2}}^{2} + \sum_{i=1}^{n} \alpha_{i}$$

$$-\frac{1}{2}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_{i}\alpha_{j}Y_{i}Y_{j}\cdot k(X_{i},X_{j})$$

$$=-\frac{1}{2\lambda}\sum_{i=1}^{n}\sum_{j=1}^{n}\alpha_{i}\alpha_{j}\cdot\left(\gamma_{i}\gamma_{j}\cdot k(x_{i},x_{j})\right)-\frac{1}{2}\|\alpha+\beta\|_{L_{2}}^{2}+1^{T}\alpha$$

$$= -\frac{1}{2\lambda} \alpha^{T} Q \alpha - \frac{1}{2} \|\alpha + \beta\|_{L_{2}}^{2} + 1 \alpha$$

Le problème dual à (Pi) est

maximiser
$$-\frac{1}{2\pi} \alpha^{T} Q \alpha - \frac{1}{2} \|\alpha + \beta\|_{L_{2}}^{2} + 1 \alpha$$
 $\alpha \in \mathbb{R}^{2}$, $\beta \in \mathbb{R}^{2}$

Pour
$$y(\beta) = -\frac{1}{2} \|\alpha + \beta\|_{\ell_{2}}^{2}$$
, $\nabla_{\beta} y = -\beta^{T}(\alpha + \beta)$

$$\nabla_{\beta} \mathcal{G} = 0 \implies \beta^* = -\alpha \leq 0$$

comme 9(B) est concave, on a

日期:

 $\max \mathcal{Y}(u) = \mathcal{Y}(0) = -\frac{1}{2} \|\alpha\|_{\ell_{x}}^{2}$

Donc le problème dual est équivalent à maximiser $-\frac{1}{2\hbar} \alpha^{2}Q\alpha + 1 \alpha - \frac{1}{2} \alpha^{2}\alpha$ $\alpha \in \mathbb{R}^{n}_{+}$

 \iff

minimiser $\frac{1}{23}\alpha^{T}(Q+3.1n)\alpha - 1\alpha$

avec P = Q+7.In

3. (1 point) Montrer que P est symétrique et semi-définie positive. Que peut-on en déduire de (P2)?

(3) D'après (2), P=Q+>In avec Qij=YiYj·k(Xi,Xj)

Symétrique

comme $\lambda \ln$ est symétrique, Q est symétrique on a $P = Q + \lambda \ln$ est symétrique

Semi-définie positive

comme $\lambda > 0$, on a $\lambda = 1$ est semi-définie positive d'après la définition $K = (k(X_i, X_j))_{1 \le j \le n}$ est S.D.P et comme $\forall i,j$ $X_i Y_j = Y_j Y_i = \pm 1$

on a Q est semi-définie positive donc P=Q+>In est semi-définie positive

	++-0	
_	нн	
	ДН	
_	77/]	

Pour (P_{2}) , comme P est semi-définie positive, alors $\frac{1}{2}\alpha^{T}P\alpha-1|^{T}\alpha$ est convexe

donc il existe un minimiseur pour $\frac{1}{7}\alpha^T P\alpha - 1^T\alpha$ $\nabla_{\alpha} \left(\frac{1}{7}\alpha^T P\alpha - 1^T\alpha\right) = P\alpha - 1$

on a $\alpha^* = P^{-1}1$ est une solution de (P_2)

4. (1 point) Expliciter les étapes d'un algorithme de résolution de (P2) de type « descente par coordonnée ».

(4) Algorithm: Sequential minimal optimization

Input: C>0, $k: \mathbb{R}^n \times \mathbb{R}^n \longrightarrow \mathbb{R}$, $\{(X_i,Y_i)\}_{1 \leq i \leq n}$

Q \leftarrow $(\gamma_i \gamma_j \cdot k(x_i, x_j))_{i \leq i, j \leq n}$

while not converge do:

find a for which KKT condition are violated

pick aj * ai at random

solve (Pz) with respect to (Qi, Qj) with all other fixed end while

Output: (a,..., an)

(5) <u>Les conditions KKT</u>
Si (P1) est convexe et les conditions de qualification
de Slater sont vérifiées.
Abors (h*, 3*) est solution de 1P1) et (a*, p*) est
solution de (Pz), ssi
$\mathbb{C} \forall i \in [1,n], g_i(h, 3) \leq 0, m_i(4) \leq 0$
② α* ≥ 0
3 $\forall i \in [1,n]$, $\alpha_i = 0$ on $g_i(h, 3) = 0$
@ Vr L(1,3, a, B) = 0 Vz L(1,3,a,B) = 0
Donc on a $h^* = \frac{1}{\lambda} \sum_{i=1}^{n} \alpha_i^* (Y_i \cdot k(\cdot, X_i))$
alors la classifieur est
$\hat{g}^*(x) = sign(h^*(x))$

5. (1½ points) Énoncer les conditions KKT pour des candidats solutions (h^*, ξ^*) et (α^*, β^*) et en déduire une classifieur issu de la résolution de (P1).

(b) Si on connaît h^* , d'après (P_1) on a $\forall i \in [1,n]$ $\begin{cases} 3_i > 1 - \gamma_i h(x_i) \\ 3_i > 0 \end{cases}$	
donc on peut choisir $\frac{3}{3}i = \max(0, 1 - \frac{1}{2}ih^*(xi))$, $\forall i \in \mathbb{C}$	زِير
$\forall i \in \{i,n\}$, si $\forall i \wedge h^*(x_i) > 1$, alors $\vec{z}_i^* = 0$ et comme $\vec{z}_i^* = \alpha_i^* + \beta_i^*$ et $\alpha_i^* > 0$, $\beta_i^* = 0$ alors $\alpha_i^* = 0$	
7. (1 point (bonus)) Proposer un critère d'arrêt pour l'algorithme itératif de la qu tion 4. Le justifier rapidement et expliciter le calcul.	es-
17) le critère d'arrêt est que tous les où sont vérifiées les conditions de KKT.	

6. (1½ points) Justifier que, connaissant h^* , on peut choisir $\xi_i^* = \max(0, 1 - Y_i h^*(x_i))$ pour tout $i \in [1, n]$. En déduire que pour tout $i \in [1, n]$, si $Y_i h^*(x_i) > 1$, alors

 $\alpha_i^{\star} = 0.$

日期:		