

EXAMEN PARTIEL

Instructions : – Une feuille aide-mémoire recto verso manuscrite est permise ;

– Durée de l'examen : 2 h 50.

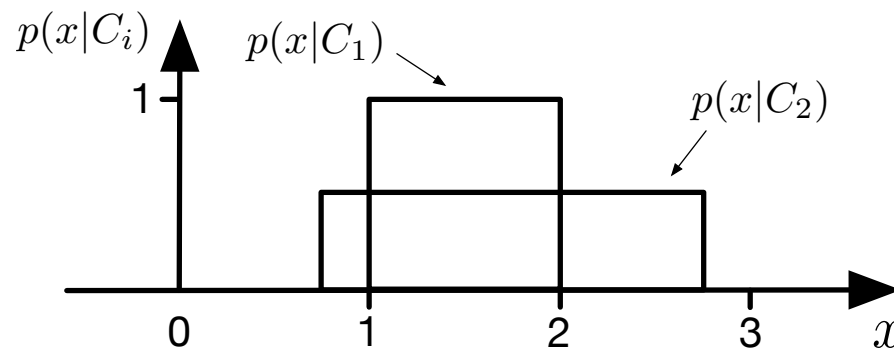
Pondération : Cet examen compte pour 30% de la note finale.

Question 1 (15 points sur 100)

Soit un problème de classement à deux classes et en une dimension, dont les vraisemblances de classe sont décrites par les densités de probabilité suivantes :

$$p(x|C_1) = \begin{cases} 1 & x \in [1, 2[\\ 0 & \text{autrement} \end{cases}, \quad p(x|C_2) = \begin{cases} 0,5 & x \in [0,75, 2,75[\\ 0 & \text{autrement} \end{cases},$$

qui sont représentées dans ce qui suit.



- (5) (a) Supposons que les probabilités a priori sont égales ($P(C_1) = P(C_2) = 0,5$), donnez la règle de décision à utiliser selon cette modélisation pour assigner une donnée x à la classe C_1 ou C_2 .

Solution:

$$h(x) = \begin{cases} C_1 & x \in [1, 2[\\ C_2 & x \in [0,75, 1[\text{ ou } x \in [2, 2,75[\\ \text{indéfini} & \text{autrement} \end{cases} \quad (1)$$

- (5) (b) Toujours en supposant des probabilités a priori égales, donnez les équations des probabilités a posteriori ($P(C_i|x)$) pour les deux classes.

Solution:

$$P(C_1|x) = \begin{cases} 0,667 & x \in [1, 2[\\ 0 & \text{autrement} \end{cases}$$

$$P(C_2|x) = \begin{cases} 1 & x \in [0,75, 1[\text{ ou } x \in [2, 2,75[\\ 0,333 & x \in [1, 2[\\ 0 & \text{autrement} \end{cases}$$

- (5) (c) Supposons maintenant que les probabilités a priori sont différentes, avec $P(C_1) = 0,25$ et $P(C_2) = 0,75$. Donnez les probabilités a posteriori pour les deux classes et la règle de décision associée.

Solution:

$$P(C_1|x) = \begin{cases} 0,4 & x \in [1, 2[\\ 0 & \text{autrement} \end{cases}$$

$$P(C_2|x) = \begin{cases} 1 & x \in [0,75, 1[\text{ ou } x \in [2, 2,75[\\ 0,6 & x \in [1, 2[\\ 0 & \text{autrement} \end{cases}$$

$$h(x) = \begin{cases} C_2 & x \in [0,75, 2,75[\\ \text{indéfini} & \text{autrement} \end{cases}$$

Question 2 (25 points sur 100)

Soit un réseau de neurones de type RBF pour deux classes, composé d'une couche cachée de R neurones de type gaussien, suivi d'une couche de sortie d'un neurone avec fonction d'activation linéaire. La valeur de la sortie pour un tel réseau de neurones pour une valeur d'entrée \mathbf{x} est donnée par l'équation suivante,

$$h(\mathbf{x}) = \sum_{i=1}^R w_i \phi_i(\mathbf{x}) + w_0 = \sum_{i=1}^R w_i \exp \left[-\frac{\|\mathbf{x} - \mathbf{m}_i\|^2}{2s_i^2} \right] + w_0,$$

où :

- \mathbf{m}_i est la valeur du centre du i -ème neurone gaussien de la couche cachée ;
- s_i est l'étalement du i -ème neurone gaussien ;
- w_i est le poids connectant le i -ème neurone gaussien de la couche cachée au neurone de sortie ;

— w_0 est le biais du neurone de sortie.

Supposons que l'on fixe les étalements s_i à des valeurs prédéterminées et que l'on veut apprendre les valeurs w_i , w_0 et m_i par descente du gradient, en utilisant comme critère l'erreur quadratique moyenne,

$$E = \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} (e^t)^2 = \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} [r^t - h(\mathbf{x}^t)]^2,$$

où :

- $r^t \in \mathbb{R}$ est la valeur désirée pour le neurone de sortie du réseau ;
- \mathcal{X} est l'ensemble des N données d'entraînement.

- (13) (a) Développez les équations permettant de **mettre à jour les poids** w_i et w_0 du neurone de sortie par descente du gradient, avec un taux d'apprentissage η , en utilisant le critère de l'erreur quadratique moyenne.

Solution:

$$\begin{aligned} e^t &= r^t - h(\mathbf{x}^t) = r^t - \left[\sum_{j=1}^R w_j \phi_j(\mathbf{x}^t) + w_0 \right] \\ \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} (e^t)^2 = \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} 2e^t \frac{\partial}{\partial w_i} \left(r^t - \left[\sum_{j=1}^R w_j \phi_j(\mathbf{x}^t) + w_0 \right] \right) \\ &= -\frac{1}{N} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t \phi_i(\mathbf{x}^t) \\ \Delta w_i &= -\eta \frac{\partial E}{\partial w_i} = \frac{\eta}{N} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t \phi_i(\mathbf{x}^t) \\ \frac{\partial E}{\partial w_0} &= \frac{\partial}{\partial w_0} \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} (e^t)^2 = \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} 2e^t \frac{\partial}{\partial w_0} \left(r^t - \left[\sum_{j=1}^R w_j \phi_j(\mathbf{x}^t) + w_0 \right] \right) \\ &= -\frac{1}{N} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t \\ \Delta w_0 &= -\eta \frac{\partial E}{\partial w_0} = \frac{\eta}{N} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t \\ w_i &= w_i + \Delta w_i, \quad i = 0, \dots, R \end{aligned}$$

- (12) (b) Développez les équations permettant de **mettre à jour les valeurs des centres** \mathbf{m}_i des neurones gaussiens de la couche cachée par descente du gradient, en utilisant le critère de l'erreur quadratique moyenne.

Solution:

$$\begin{aligned}
 \frac{\partial \phi_i(\mathbf{x}^t)}{\partial m_{i,j}} &= \frac{\partial}{\partial m_{i,j}} \exp \left[-\frac{\|\mathbf{x}^t - \mathbf{m}_i\|^2}{2s_i^2} \right] \\
 &= \exp \left[-\frac{\|\mathbf{x}^t - \mathbf{m}_i\|^2}{2s_i^2} \right] \frac{\partial}{\partial m_{i,j}} \left[-\frac{\|\mathbf{x}^t - \mathbf{m}_i\|^2}{2s_i^2} \right] \\
 &= \frac{(x_j^t - m_{i,j})}{s_i^2} \exp \left[-\frac{\|\mathbf{x}^t - \mathbf{m}_i\|^2}{2s_i^2} \right] = \frac{x_j^t - m_{i,j}}{s_i^2} \phi_i(\mathbf{x}^t) \\
 \frac{\partial E}{\partial m_{i,j}} &= \frac{\partial}{\partial m_{i,j}} \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} (e^t)^2 = \frac{1}{2N} \sum_{\mathbf{x}^t \in \mathcal{X}} 2e^t \frac{\partial}{\partial m_{i,j}} \left(r^t - \left[\sum_{l=1}^R w_l \phi_l(\mathbf{x}^t) + w_0 \right] \right) \\
 &= \frac{1}{N} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t (-1) w_i \frac{\partial \phi_i(\mathbf{x}^t)}{\partial m_{i,j}} = -\frac{w_i}{N s_i^2} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t (x_j^t - m_{i,j}) \phi_i(\mathbf{x}^t) \\
 \Delta m_{i,j} &= -\eta \frac{\partial E}{\partial m_{i,j}} = \frac{\eta w_i}{N s_i^2} \sum_{\mathbf{x}^t \in \mathcal{X}} e^t (x_j^t - m_{i,j}) \phi_i(\mathbf{x}^t) \\
 m_{i,j} &= m_{i,j} + \Delta m_{i,j}, \quad i = 1, \dots, R, \quad j = 1, \dots, D
 \end{aligned}$$

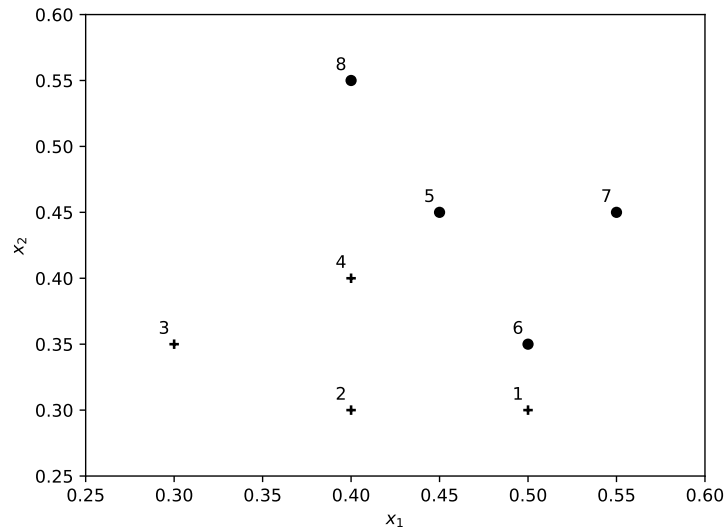
Question 3 (32 points sur 100)

Soit le jeu de données suivant, en deux dimensions :

$$\begin{aligned}
 \mathbf{x}^1 &= [0,5 \ 0,3]^\top, & \mathbf{x}^2 &= [0,4 \ 0,3]^\top, & \mathbf{x}^3 &= [0,3 \ 0,35]^\top, & \mathbf{x}^4 &= [0,4 \ 0,4]^\top, \\
 \mathbf{x}^5 &= [0,45 \ 0,45]^\top, & \mathbf{x}^6 &= [0,5 \ 0,35]^\top, & \mathbf{x}^7 &= [0,55 \ 0,45]^\top, & \mathbf{x}^8 &= [0,4 \ 0,55]^\top.
 \end{aligned}$$

Les étiquettes de ces données sont $r^1 = r^2 = r^3 = r^4 = -1$ et $r^5 = r^6 = r^7 = r^8 = 1$.

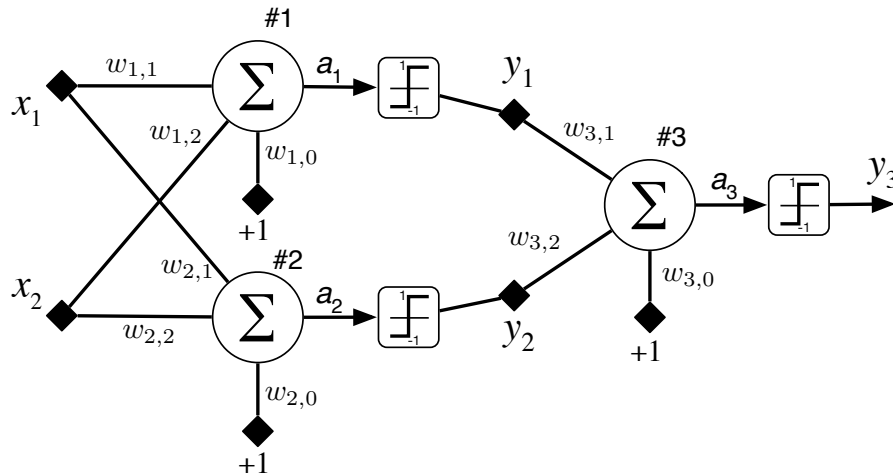
Le graphique ici-bas présente le tracé de ces données.



- (10) (a) Soit un réseau de neurones utilisant la fonction signe comme fonction d'activation :

$$f_{\text{sgn}}(a) = \text{sgn}(a) = \begin{cases} 1 & a \geq 0 \\ -1 & a < 0 \end{cases}.$$

Supposons que l'on veut entraîner le réseau suivant à trois neurones avec cette fonction d'activation signe, pour classifier les données présentées ci-haut, selon le schéma suivant.



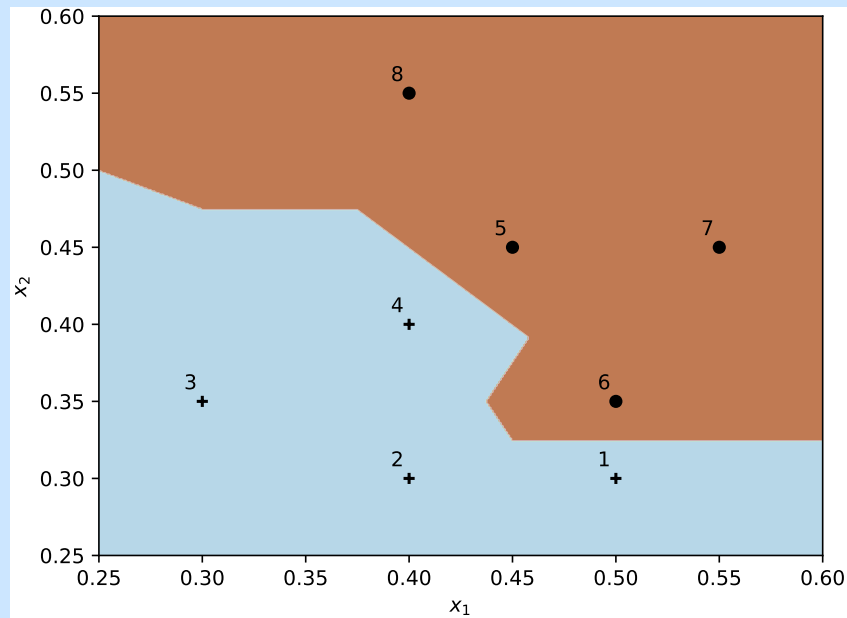
Supposons les valeurs de poids $w_{3,1} = 1$ et $w_{3,2} = 1$, et du biais $w_{3,0} = -1,5$. Déterminez les poids et biais des neurones 1 et 2 afin d'obtenir un taux de classement correct de 100 % sur les données.

Solution: Plusieurs solutions sont possibles. Une configuration valide est :

$$w_{1,1} = 1, \quad w_{1,2} = 1, \quad w_{1,0} = -0.825, \quad w_{2,1} = 0, \quad w_{2,2} = 1, \quad w_{2,0} = -0.325.$$

- (10) (b) Tracez les régions de décision selon les données d'entraînement présentée en préambule de la question pour un classifieur de type plus proche voisin (avec un seul voisin, $k = 1$) utilisant une distance euclidienne. Tracez le tout dans la **feuille de réponse fournie** (et non dans l'énoncé courant). Donnez également le taux de classement selon une méthodologie *leave-one-out* avec cette configuration, sur ces données.

Solution:



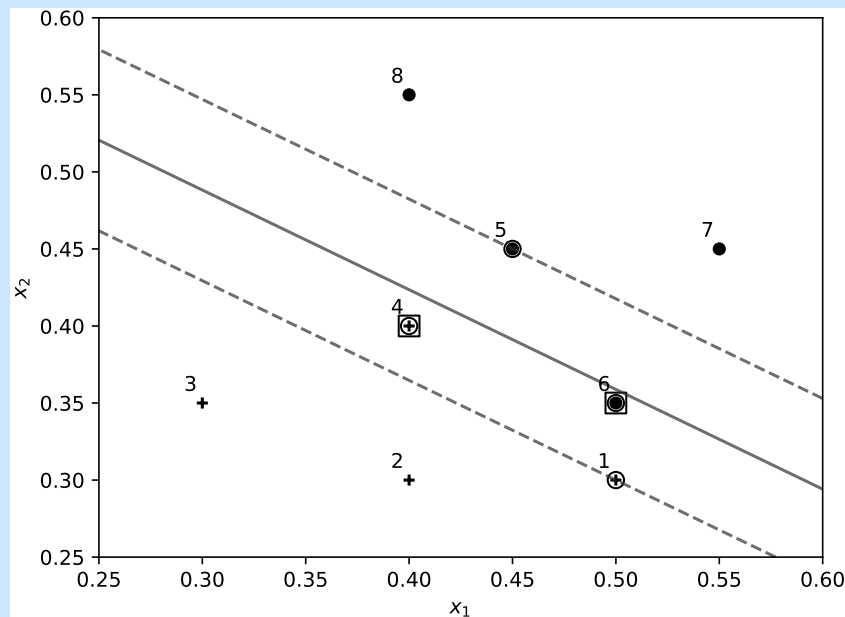
Données mal classées par *leave-one-out* : x^1, x^4, x^5, x^6 . Taux de classement correspondant de 50 % (4/8).

- (12) (c) Nous obtenons le résultat suivant en effectuant l'entraînement d'un SVM linéaire à **marge douce** avec ces données, avec comme valeur de paramètre de régularisation $C = 200$:

$$\alpha^1 = 180, \quad \alpha^2 = 0, \quad \alpha^3 = 0, \quad \alpha^4 = 200, \quad \alpha^5 = 180, \quad \alpha^6 = 200, \quad \alpha^7 = 0, \quad \alpha^8 = 0, \\ w_0 = -11,6.$$

En utilisant la feuille de réponse fournie, tracez les éléments suivants :

- Frontière de décision du SVM (droite continue, —);
- Frontières de la marge (droites en traits pointillés, ----);
- Vecteurs de support (encerclez les points, ○);
- Données dans la marge (encadrez les points, □).

Solution:

Question 4 (28 points sur 100)

Répondez aussi brièvement et clairement que possible aux questions suivantes.

- (4) (a) Expliquez en quoi consiste la régularisation dans un contexte d'apprentissage supervisé.

Solution: La régularisation consiste à inclure une mesure sur la complexité du modèle dans le processus d'apprentissage, afin qu'à performance similaire, des modèles plus simples soient favorisés.

- (4) (b) Dans un contexte de classement paramétrique s'appuyant sur des modèles de densité de probabilité des données selon les classes, la règle de Bayes est souvent utilisée de la façon suivante :

$$\underbrace{P(C|\mathbf{x})}_{\text{a posteriori}} = \frac{\overbrace{P(C)}^{\text{a priori}} \overbrace{p(\mathbf{x}|C)}^{\text{vraisemblance}}}{\underbrace{p(\mathbf{x})}_{\text{évidence}}}.$$

Expliquez pourquoi l'évidence $p(\mathbf{x})$ de cette formulation est généralement ignorée lorsqu'il vient le moment de définir la règle de décision (fonction $h_i(\mathbf{x})$).

Solution: L'évidence est la même pour toutes les classes alors qu'on s'intéresse qu'à la classe pour laquelle la fonction de décision $h_i(\mathbf{x})$ est maximale. Donc, peu importe la valeur de l'évidence $p(\mathbf{x})$, qui est la même pour toutes les classes, les décisions basées sur la valeur maximale de fonctions $h_i(\mathbf{x}) = P(C_i) p(\mathbf{x}|C_i)$ seront les mêmes qu'une décision basée sur la valeur maximale de fonctions $h_i(\mathbf{x}) = P(C_i|\mathbf{x})$.

- (4) (c) Supposons un classement paramétrique basé sur des lois normales multivariées, où l'estimation de la matrice de covariance \mathbf{S} est partagée (la même) entre toutes les classes. Indiquez l'équation pour calculer cette estimation de la matrice de covariance partagée à partir d'un jeu de données étiquetées pour le classement.

Solution:

$$\mathbf{S} = \sum_{i=1}^K P(C_i) \mathbf{S}_i,$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^\top}{\sum_t r_i^t}.$$

- (4) (d) Dans un contexte de classement par les k -plus proches voisins, expliquez l'effet du nombre de voisins k sur le classement et les frontières de décision.

Solution: Lorsque k est faible, les frontières de décisions peuvent être plus bruitées, mais peuvent également capturer des éléments présents dans plus de données. Avec un k plus élevé, les frontières sont plus fortement adoucies, réduisant la sensibilité aux données aberrantes, mais également réduisant la capacité à modéliser des éléments significatifs présents dans seulement quelques instances de données.

- (4) (e) Expliquez le lien principal que l'on peut établir entre la régression logistique et le classement paramétrique.

Solution: La régression logistique fait une approximation par un modèle linéaire joint à une fonction sigmoïde de la probabilité a posteriori $P(C_i|\mathbf{x})$, qui peut s'interpréter de façon similaire à du classement paramétrique avec des lois normales, où les matrices de covariances sont partagées entre les classes.

- (4) (f) Dans un contexte de classement avec SVM à marges douces, expliquez comment on interprète les variables *slacks* ξ^t .

Solution: Les variables $\xi^t > 0$ dénotent des données qui sont dans la marge du SVM, c'est-à-dire des données pour lequel on obtient que $r^t h(\mathbf{x}^t) < 1$. Pour $0 < \xi^t < 1$, la donnée est bien classée mais dans la marge, alors que pour $\xi^t \geq 1$, la donnée est mal classée.

- (4) (g) Expliquez la différence entre un apprentissage par lots (*batch*) et un apprentissage en ligne dans un réseau de neurones.

Solution: Un apprentissage par lot offre une stabilité à l'apprentissage en calculant des corrections à appliquer aux poids et biais du réseau de neurones sur l'ensemble du jeu de données d'entraînement. Cependant, cet apprentissage peut prendre plus d'époques à converger qu'un apprentissage en ligne, où les corrections sont calculées pour chaque donnée du jeu d'entraînement.