

Introduction à l'apprentissage automatique

Tronc commun scientifique 2A - Mines Nancy

Examen - mardi 25 janvier 2022

Durée : 2 heures 30 minutes

Ce test comporte deux parties. À la fin du test, rendez votre copie (partie 1) au surveillant et déposez sur Arche votre carnet Jupyter (partie 2).

La qualité et la précision de la rédaction seront valorisées.

Note finale au module :

$$\text{Note de TD (5 points)} + \frac{15}{20} (\text{Note partie 1 (10 points)} + \text{Note partie 2 (10 points)})$$

Partie 1

Justifiez brièvement, sauf pour la question 1 dans laquelle vous donnerez seulement la réponse.

1) 1pt Indiquez pour chaque méthode d'apprentissage si elle est supervisée ou non-supervisée :

- (a) k plus proches voisins
- (b) machines à vecteurs support
- (c) perceptrons multi-couches
- (d) classification hiérarchique ascendante

2) 1,5pt

- (a) dans le cas biclasse, un classifieur linéaire sépare toujours les observations des deux classes par un hyperplan.

vrai / faux

- (b) les SVM fournissent des probabilités d'appartenance aux classes, comme le classifieur de la régression logistique.

vrai / faux

- (c) les SVM linéaires fournissent l'erreur de généralisation la plus faible parmi les classifieurs linéaires.

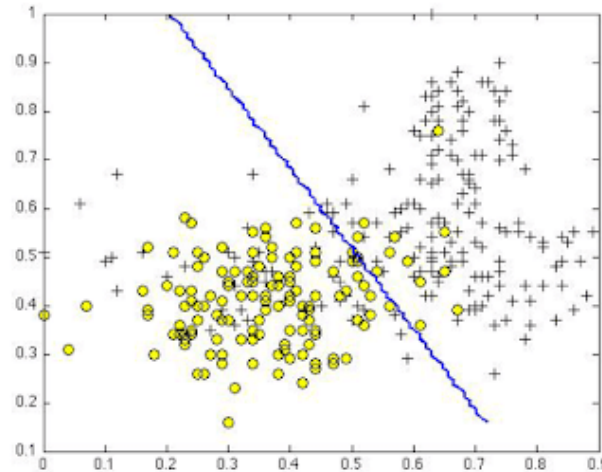
vrai / faux

3) 0,5pt On considère un problème de classification biclasse (YES ou NO) pour lequel un modèle d'apprentissage fournit la matrice de confusion suivante sur la base test :

n=165	Predicted: NO	Predicted: YES
Actual: NO	50	10
Actual: YES	5	100

Quel est le taux de classifications correctes du modèle? (*accuracy score*)

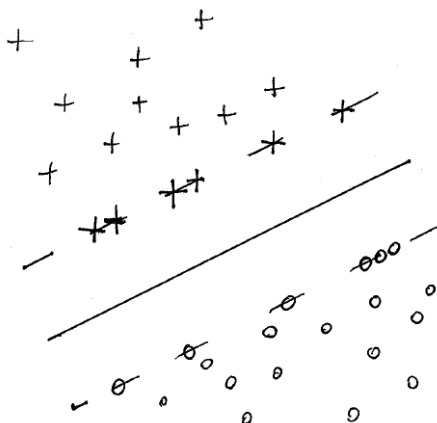
- 4) 1pt On dispose de la matrice D des distances entre cent observations mais pas des observations elles-mêmes. Autrement dit, l'élément de ligne i et colonne j de la matrice D est la distance entre les observations d'indices i et j . Quel(s) algorithme(s) de partitionnement vu(s) en cours ou TD pouvez-vous mettre en œuvre, et pourquoi?
- 5) 1pt On considère un problème de classification à deux classes (o et +) d'observations bidimensionnelles. Une classification par SVM à noyau RBF (gaussien, de paramètre γ) fournit la « courbe » de séparation des classes suivante :



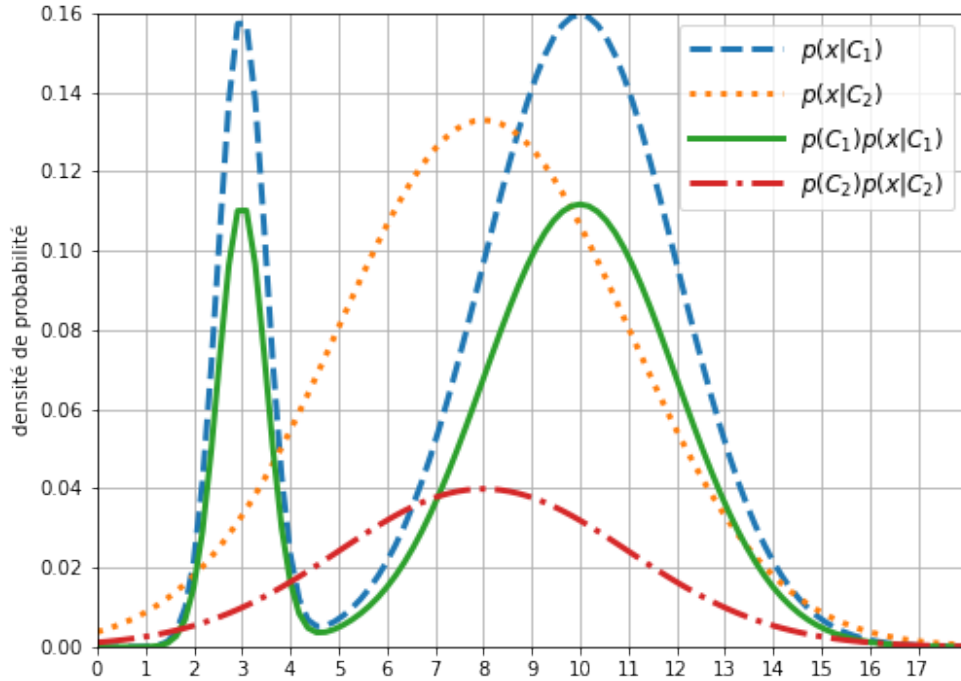
Est-on dans une situation de sous-apprentissage ou de surapprentissage? Que pouvez-vous essayer pour améliorer les performances de ce modèle sans modifier la valeur de l'hyperparamètre γ ?

- 6) 1pt On considère le problème de classification supervisée représenté par la figure suivante, où on peut voir les observations d'apprentissage (classes + et o) et la droite séparatrice de marge maximale obtenue sur l'intégralité de l'ensemble d'apprentissage. Quel est le score de validation croisée *leave-one-out* du modèle des séparateurs à vaste marge (SVM à noyau linéaire sans variables d'écart)? On utilise comme score le taux de classifications correctes.

Rappel: la validation *leave-one-out* correspond à la validation croisée à N plis, où N est le nombre d'observations dans la base d'apprentissage.



- 7) 2pt On considère un problème de classification biclasse sur \mathbb{R} pour lequel on connaît différentes distributions de probabilité représentées sur le graphique suivant, ainsi que les probabilités a priori $p(C_1) = 0,7$ et $p(C_2) = 0,3$.



On souhaite prédire la classe d'une observation $x \in [0, 16]$. Quelle règle de classification correspond au classifieur de Bayes? (vous vous conterez de lire des valeurs approchées sur le graphique)

- 8) 2pts On considère l'ensemble d'observations unidimensionnelles $D = \{1, 2, 3, 5, 8, 9, 13, 14\}$. On cherche un partitionnement de D en trois groupes, à l'aide de la méthode des K -moyennes (K -means). La distance entre deux observations a et b réelles est simplement $d(a, b) = |a - b|$. Détaillez les différentes étapes de l'algorithme de Lloyd en choisissant comme centres initiaux des groupes les observations $\mu_1 = 1$ (groupe G_1), $\mu_2 = 8$ (G_2), et $\mu_3 = 9$ (G_3), et donnez la classification obtenue.

Donnez une autre initialisation qui aurait fourni un partitionnement à trois classes différent.

À l'aide de quel critère pourriez-vous choisir l'une ou l'autre des partitions? (on ne vous demande pas de calcul, seulement d'expliquer le principe de choix)

- 9) question subsidiaire, bonus de 2pts On considère un problème de régression linéaire univarié : chacune des N observations est une donnée réelle x_i , et l'étiquette associée est un réel y_i . Le coût des erreurs est la somme des carrés des résidus du modèle. En quoi consiste l'algorithme du gradient stochastique dans ce cas?

(remarque : la démarche se généralise au cas multivarié, et est utile quand la taille du jeu de données rend difficile l'inversion de la matrice qui intervient dans les équations normales)

Partie 2

10pts Étude du jeu de données `glass2.txt`, sur Arche.

Déposez votre carnet avec les cellules exécutées au format HTML (File → Download as → HTML). On doit pouvoir vous évaluer à partir de ce fichier HTML visualisé dans un navigateur. Préparez votre dépôt quelques minutes avant la fin du test.