

Consignes :

- Sont interdits : Documents, calculatrices, téléphones, écouteurs, ordinateurs, tablettes.
- Il est interdit de composer avec un crayon.
- Votre feuille double d'examen doit porter, à l'emplacement réservé, vos nom, prénom, et signature.
- Cette zone réservée doit être cachée par collage.
- Vos feuilles intercalaires doivent être toutes numérotées.
- Le barème est donné à titre indicatif.

Exercice 1 (4 pts)

1. Présenter en trois lignes le principe d'une analyse prédictive de données et celui d'une analyse descriptive.
2. Présenter en six lignes une méthode d'apprentissage statistique supervisée et une non-supervisée : Pour cela on décrira le principe général en parlant du fondement probabiliste, du paradigme génératif ou discriminatif, de la nature du critère optimisé (dans le cas où il y en a un), de l'existence ou non d'une solution analytique, de sa mise en œuvre algorithmique, de la notion de convergence locale ou globale, etc)

Exercice 2 (9 pts) On considère un échantillon indépendant $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ d'individus \mathbf{X}_i décrits par d variables binaires $(\mathbf{X}_i = (X_{i1}, \dots, X_{id})^T \in \{0, 1\}^d)$ indépendantes. On suppose que ces individus sont issus d'une population hétérogène à K classes inconnues. Soit $Y_i \in \llbracket 1, K \rrbracket$ la classe inconnue de l'individu \mathbf{X}_i . On dispose d'un échantillon observé $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ où $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^T \in \{0, 1\}^d$ est la réalisation de \mathbf{X}_i . On suppose que les individus $(\mathbf{X}_1, \dots, \mathbf{X}_n)$ sont identiquement distribués selon la loi mélange de lois de Bernoulli multivariées définie par :

$$\mathbb{P}(\mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{B}(\mathbf{x}_i; \mathbf{p}_k) \quad (1)$$

où $\pi_k > 0$ sont les proportions du mélange vérifiant $\sum_{k=1}^K \pi_k = 1$, et $\mathcal{B}(\mathbf{x}_i; \mathbf{p}_k)$ la loi de la k ème classe.

1. Justifier pourquoi a-t-on $\mathcal{B}(\mathbf{x}_i; \mathbf{p}_k) = \prod_{j=1}^d \mathcal{B}(x_{ij}; p_{kj})$ où $\mathcal{B}(x_{ij}; p_{kj})$ est la loi de Bernoulli univariée de paramètre $p_{kj} \in]0, 1[$ associée à la variable x_{ij} et définie par $\mathcal{B}(x; p) = p^x (1 - p)^{1-x}$.
2. Définir le vecteur paramètre $\boldsymbol{\theta}$ et donner le nombre de paramètres libres $\nu_{\boldsymbol{\theta}}$ du modèle
3. On note par $L(\boldsymbol{\theta})$ (respectivement $\log L(\boldsymbol{\theta})$) la vraisemblance (respectivement log-vraisemblance) de $\boldsymbol{\theta}$ pour les données observées $(\mathbf{x}_1, \dots, \mathbf{x}_n)$. Montrer que la maximisation de la log-vraisemblance $\log L(\boldsymbol{\theta})$ par l'algorithme EM conduit aux étapes suivantes à chaque itération (t) :

I.Étape E : Cette étape nécessite seulement le calcul des probabilités suivantes :

$$\tau_{ik}^{(t)} = \frac{\pi_k^{(t)} \mathcal{B}(\mathbf{x}_i; \mathbf{p}_k^{(t)})}{\sum_{\ell=1}^K \pi_{\ell}^{(t)} \mathcal{B}(\mathbf{x}_i; \mathbf{p}_{\ell}^{(t)})}$$

II.Étape M : La mise à jour des paramètres du modèle $\boldsymbol{\theta}^{(t+1)}$ s'effectue selon les formules suivantes :

$$\pi_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(t)}}{n}$$

et

$$\mathbf{p}_k^{(t+1)} = \frac{\sum_{i=1}^n \tau_{ik}^{(t)} \mathbf{x}_i}{\sum_{i=1}^n \tau_{ik}^{(t)}}$$

4. Après la convergence de l'algorithme, on obtient le vecteur paramètre estimé du modèle $\hat{\boldsymbol{\theta}}$. En déduire l'expression d'une règle pour déterminer une estimation de la classe \hat{y}_i d'un individu \mathbf{x}_i .
5. On dispose des modèles estimés décrits par $\hat{\boldsymbol{\theta}}_1, \dots, \hat{\boldsymbol{\theta}}_{Kmax}$ où $\hat{\boldsymbol{\theta}}_K, K \in \llbracket 1, Kmax \rrbracket$, correspond à l'estimation par maximum de vraisemblance du vecteur paramètre d'un modèle de la forme de (1). Donner l'expression d'un critère pour sélectionner le nombre optimal des classes \hat{K} .

Exercice 3 (7 pts)

0.5 pt si réponse correcte, -0.5 pt si réponse incorrecte, 0 en cas de non réponse.

Une seule réponse est correcte.

Feuille à remettre. NE PAS INSCRIRE VOS NOM ET PRENOM SUR CETTE FEUILLE

Questions	Réponses
1. K -nn (K -ppv) est une méthode de classification par apprentissage	<input type="checkbox"/> Oui <input type="checkbox"/> Non
2. K -nn (K -ppv) est un classifieur qui passe facilement à l'échelle	<input type="checkbox"/> Oui <input type="checkbox"/> Non
3. La régression logistique est une méthode	<input type="checkbox"/> supervisée <input type="checkbox"/> non supervisée
4.	<input type="checkbox"/> générative <input type="checkbox"/> discriminative
5.	<input type="checkbox"/> qui admet une solution analytique <input type="checkbox"/> qui nécessite un algorithme d'optimisation
6. L'analyse linéaire discriminante (LDA) est une méthode	<input type="checkbox"/> supervisée <input type="checkbox"/> non supervisée
7.	<input type="checkbox"/> générative <input type="checkbox"/> discriminative
8. En LDA, l'estimation des paramètres	<input type="checkbox"/> nécessite un algorithme d'optimisation <input type="checkbox"/> s'effectue de façon exacte
9. En LDA, la vraisemblance maximisée est	<input type="checkbox"/> concave <input type="checkbox"/> non-concave
10. LDA suppose une matrice de covariance	<input type="checkbox"/> différente pour toutes les classes <input type="checkbox"/> commune à toutes les classes
11. K -means est un algorithme qui se base sur un formalisme probabiliste	<input type="checkbox"/> Oui <input type="checkbox"/> Non
12. K -means est insensible à l'ordre de présentation des données	<input type="checkbox"/> Oui <input type="checkbox"/> Non
13. Dans le cas de son utilisation pour un mélange Gaussien (GMM), EM est insensible à l'ordre de présentation des données	<input type="checkbox"/> Oui <input type="checkbox"/> Non
14. EM garantit l'obtention d'un	<input type="checkbox"/> optimum global <input type="checkbox"/> optimum local