

Consignes :

- Sont interdits : Documents, calequettes, 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) Soit (X, Y) un couple de variables aléatoires réelles et soit $((x_1, y_1), \dots, (x_n, y_n))$ un échantillon de n observations. Chacune des situations présentées dans la Figure 1 représente le nuage de données d'un échantillon de taille $n = 500$. Pour chaque situation, donner une valeur approchée du coefficient de corrélation linéaire empirique r et justifier votre réponse.

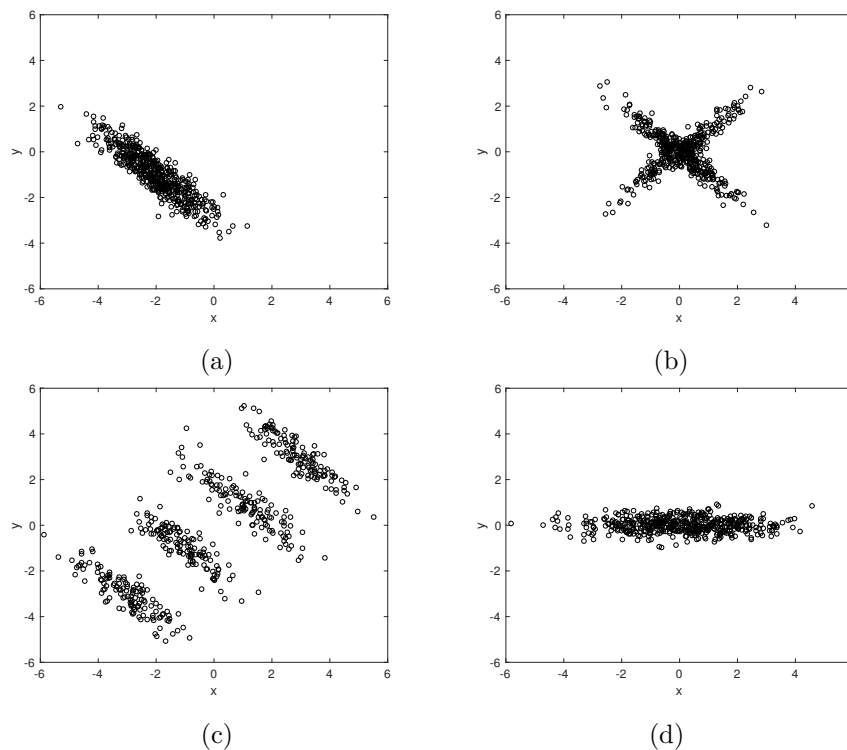


FIGURE 1 – Nuages de données (o)

Exercice 2 (4 pts) On considère un échantillon aléatoire indépendant $((\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n))$ du couple (\mathbf{X}, Y) où \mathbf{X} est un vecteur composé de $p = 2$ prédicteurs réels ($\mathbf{X} \in \mathbb{R}^2$) et $Y \in [0, 1]$ une variable à prédire. On dispose d'un échantillon observé d'apprentissage $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ et on cherche à prédire les classes de nouvelles observations $(\mathbf{x}_{n+1}, \dots, \mathbf{x}_m)$ sur la base d'un modèle probabiliste appris sur les données d'apprentissage. On considère l'analyse discriminante où la densité de la classe $k \in [0, 1]$ est définie par : $f(\mathbf{x}_i | Y_i = k; \boldsymbol{\theta}) = \frac{1}{2\pi|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k)\right)$. La prédiction s'effectue par la règle du maximum a posteriori (MAP) qui consiste à affecter l'individu \mathbf{x}_i à la classe y_i maximisant la probabilité a posteriori :

$$\hat{y}_i = \arg \max_{k \in [0,1]} \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\theta}),$$

$\boldsymbol{\theta}$ étant le vecteur paramètre du modèle. On note par $\pi_k = \mathbb{P}(Y_i = k)$, la probabilité a priori de la classe k . On suppose que $\pi_0 = 0.5$, $\boldsymbol{\mu}_0 = \boldsymbol{\mu}_1 = \mathbf{g}$, $\boldsymbol{\Sigma}_0 = \mathbf{I}$ et $\boldsymbol{\Sigma}_1 = \lambda \boldsymbol{\Sigma}_0$, où $\lambda > 1$ et \mathbf{I} est la matrice identité.

1. Montrer que $Y = 1$ si et seulement si $\|\mathbf{x} - \mathbf{g}\|_2^2 \geq r$, en déterminant r en fonction de λ .
2. A quoi correspond la frontière de décision dans ce cas ?

Exercice 3 (4 pts) On considère le cadre de l'exercice précédent et on se propose maintenant d'utiliser un autre modèle de prédiction, celui de régression logistique non-linéaire suivant défini par

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\phi(\mathbf{x}; \boldsymbol{\theta}))}{1 + \exp(\phi(\mathbf{x}; \boldsymbol{\theta}))} \quad (1)$$

où $\phi(\mathbf{x}; \boldsymbol{\theta})$ étant une transformation non-linéaire (polynomiale) de \mathbf{x} et $\boldsymbol{\theta}$ le vecteur paramètre du modèle. La prédiction avec ce modèle consiste à affecter l'individu \mathbf{x}_i à la classe y_i maximisant la probabilité a posteriori, i.e. :

$$\hat{y}_i = \arg \max_{k \in \{0,1\}} \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \hat{\boldsymbol{\theta}}),$$

où $\hat{\boldsymbol{\theta}}$ étant le vecteur paramètre du modèle appris en minimisant le risque quadratique régularisé suivant :

$$\ell_\lambda(\boldsymbol{\theta}) = -\ln L(\boldsymbol{\theta}) + \lambda \Omega(\boldsymbol{\theta}) \quad (2)$$

où $\ln L(\boldsymbol{\theta})$ est la fonction de log-vraisemblance classique et $\Omega(\boldsymbol{\theta})$ une fonction de pénalité sur $\boldsymbol{\theta}$ de niveau λ dont l'objectif est d'assurer un compromis entre complexité du modèle et ajustement aux données.

1. Rappeler le nom et le principe de l'algorithme d'estimation des paramètres en régression logistique.
2. Rappeler la décomposition biais-variance du risque quadratique.
3. Chacune des figures 2-(a, b, c) montre un échantillon d'apprentissage et un modèle de prédiction estimé $\hat{\boldsymbol{\theta}}$ en minimisant (2), représenté par la frontière de décision, pour la même régularisation $\Omega(\boldsymbol{\theta})$ et différentes valeurs de λ . Discuter la qualité prédictive de chacun des trois modèles.

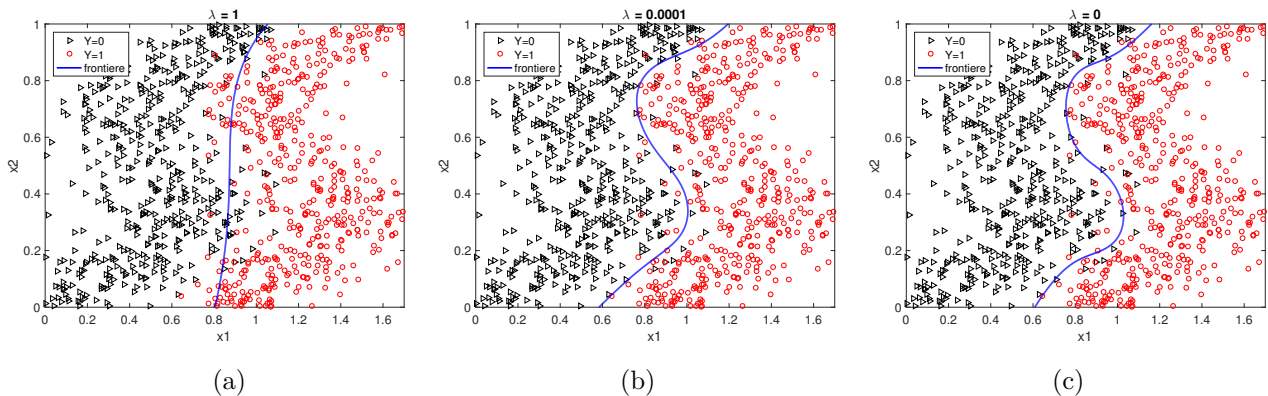


FIGURE 2 – Nuage de données (o,▷) et modèle de prédiction (—)

Exercice 4 (8 pts) On considère un échantillon aléatoire indépendant $((\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n))$ du couple (\mathbf{X}, Y) où \mathbf{X} est un vecteur de p prédicteurs binaires ($\mathbf{X} \in \llbracket 0, 1 \rrbracket^p$) et $Y \in \llbracket 0, 1 \rrbracket$ une variable à prédire représentant la classe de \mathbf{X} . On dispose d'un échantillon observé d'apprentissage $((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n))$ et on cherche à prédire les classes de nouvelles observations $(\mathbf{x}_{n+1}, \dots, \mathbf{x}_m)$ sur la base d'un modèle probabiliste appris sur les données d'apprentissage. On considère la classifieur Bayésien naïf, une méthode probabiliste de prédiction qui repose sur l'hypothèse simplificatrice forte suivante :

\mathcal{H} : Pour chaque classe Y_i , les variables $X_{ij}, j = 1, \dots, p$ de chaque individu \mathbf{X}_i , sont indépendantes.

La prédiction s'effectue par la règle du maximum a posteriori (MAP) qui consiste à affecter l'individu \mathbf{x}_i à la classe y_i maximisant la probabilité a posteriori :

$$\hat{y}_i = \arg \max_{k \in \{0,1\}} \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \hat{\Psi}), \quad (3)$$

$\hat{\Psi}$ étant l'estimateur du vecteur paramètre Ψ du modèle estimé sur l'ensemble d'apprentissage.

1. Les variables binaires X_{ij} de chaque vecteur \mathbf{X}_i de la même classe $Y_i = k$ sont de loi Bernoulli de paramètre $0 < \theta_{kj} < 1$, i.e. :

$$\mathbb{P}(X_{ij} = x_{ij} | Y_i = k; \theta_{kj}) = \theta_{kj}^{x_{ij}} (1 - \theta_{kj})^{1-x_{ij}} \text{ où } \theta_{kj} = \mathbb{P}(X_{ij} = 1 | Y_i = k; \theta_{kj}).$$

En déduire d'après \mathcal{H} la loi conditionnelle du vecteur \mathbf{X}_i , notée $\mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = k; \theta_k)$.

2. On note par $\alpha = \mathbb{P}(Y_i = 1)$, la probabilité a priori de la classe 1. On a alors $\Psi = (\alpha, \theta_1^\top, \theta_0^\top)^\top$ que l'on cherche à estimer la méthode du maximum de vraisemblance. Donner l'expression de la fonction de vraisemblance

$$L(\Psi) = \mathbb{P}((\mathbf{X}_1 = \mathbf{x}_1, Y_1 = y_1), \dots, (\mathbf{X}_n = \mathbf{x}_n, Y_n = y_n); \Psi). \quad (4)$$

3. En déduire celle de la fonction de log-vraisemblance $\ln L(\Psi)$.
4. Montrer en maximisant la log-vraisemblance que les estimateurs du maximum de vraisemblance sont donnés par :

$$\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n Y_i; \quad \hat{\theta}_{1j} = \frac{\sum_{i=1}^n Y_i X_{ij}}{\sum_{i=1}^n Y_i}; \quad \hat{\theta}_{0j} = \frac{\sum_{i=1}^n (1 - Y_i) X_{ij}}{\sum_{i=1}^n (1 - Y_i)}.$$

soit en formulation vectorielle pour les θ par :

$$\hat{\theta}_1 = \frac{\sum_{i=1}^n Y_i \mathbf{X}_i}{\sum_{i=1}^n Y_i}; \quad \hat{\theta}_0 = \frac{\sum_{i=1}^n (1 - Y_i) \mathbf{X}_i}{\sum_{i=1}^n (1 - Y_i)}.$$

5. On considère le jeu de données de la Table 1. Prédire en appliquant la règle (3) la valeur de Y pour le dernier vecteur.

\mathbf{X}	Y
0	1
1	0
1	0
0	1
0	1
1	1
1	?

TABLE 1 – Jeu de données