

**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.

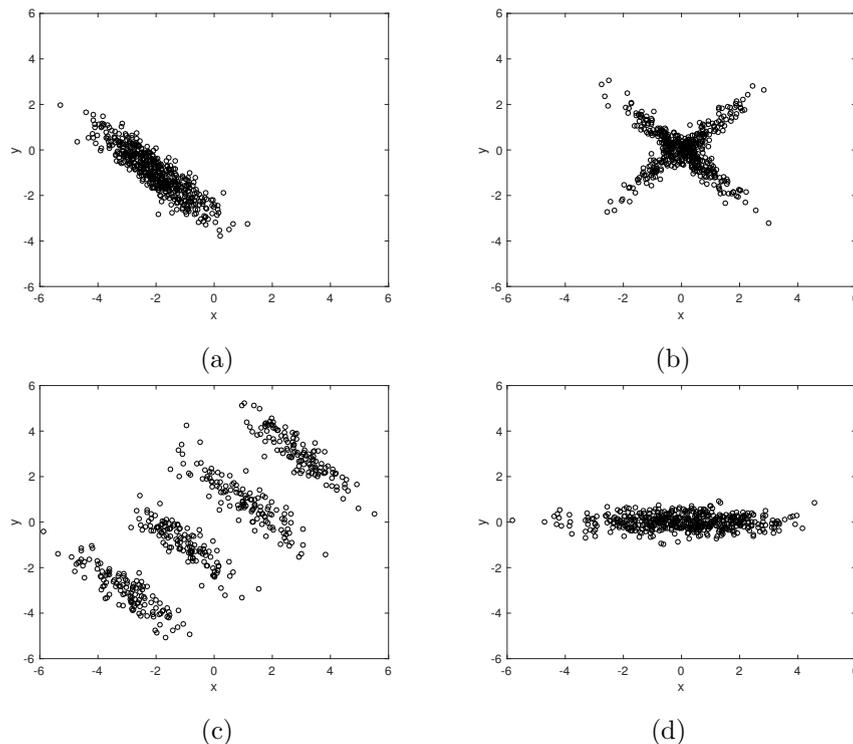


FIGURE 1 – Nuages de données (o)

**Solution 1** Voir TD ((a) forte relation linéaire (anti-corrélation),  $r$  proche de  $-0.8/0.7$ ; (b)  $r$  proche de 0 car pas de lien linéaire; (c)  $r$  proche de 0.7 (paradoxe de Simpson); (d)  $r$  proche de 0 : quasiment pas de lien entre les deux variables.

**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  $Y \in \llbracket 0, 1 \rrbracket$  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  $k \in \llbracket 0, 1 \rrbracket$  est définie par :  $f(\mathbf{x}_i | Y_i = k; \boldsymbol{\theta}) = \frac{1}{2\pi|\boldsymbol{\Sigma}_k|^{\frac{1}{2}}} \exp(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k))$ . 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 :

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

$\boldsymbol{\theta}$  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  $\lambda$ .

2. A quoi correspond la frontière de décision dans ce cas ?

**Solution 2**

1. On a  $Y = 1$  si et seulement si :

$$\begin{aligned}
 \ln \frac{\mathbb{P}(y = 1|\mathbf{x})}{\mathbb{P}(y = 0|\mathbf{x})} &= \ln \frac{\pi_1}{\pi_0} - \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_0|} - \frac{1}{2} \{(\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) - (\mathbf{x} - \boldsymbol{\mu}_0)^T \Sigma_0^{-1} (\mathbf{x} - \boldsymbol{\mu}_0)\} \geq 0 \\
 &= -\frac{1}{2} \ln(\lambda^2) - \frac{1}{2} \left\{ (\mathbf{x} - \mathbf{g})^T \frac{1}{\lambda} (\mathbf{x} - \mathbf{g}) - (\mathbf{x} - \mathbf{g})^T (\mathbf{x} - \mathbf{g}) \right\} \geq 0 \\
 &= -\ln(\lambda) - \frac{1}{2} \left\{ \|\mathbf{x} - \mathbf{g}\|_2^2 \left( \frac{1}{\lambda} - 1 \right) \right\} \geq 0 \\
 &= \frac{1}{2} \|\mathbf{x} - \mathbf{g}\|_2^2 \left( \frac{\lambda - 1}{\lambda} \right) \geq \ln(\lambda) \\
 &= \|\mathbf{x} - \mathbf{g}\|_2^2 \geq 2 \ln(\lambda) \left( \frac{\lambda}{\lambda - 1} \right)
 \end{aligned} \tag{1}$$

Au finale on a donc :  $Y = 1$  si et seulement si  $\|\mathbf{x} - \mathbf{g}\|_2^2 \geq r$  avec  $r = 2 \ln(\lambda) \left( \frac{\lambda}{\lambda - 1} \right)$ .

2. D'après la question précédent, la frontière de décision est définie par l'équation :

$$\|\mathbf{x} - \mathbf{g}\|_2^2 = 2 \ln(\lambda) \left( \frac{\lambda}{\lambda - 1} \right)$$

qui est celle d'un cercle de centre  $\mathbf{g}$  et de rayon  $\sqrt{2 \ln(\lambda) \left( \frac{\lambda}{\lambda - 1} \right)}$ .

**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

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\phi(\mathbf{x}; \boldsymbol{\theta}))}{1 + \exp(\phi(\mathbf{x}; \boldsymbol{\theta}))} \tag{2}$$

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 :

$$\ell_\lambda(\boldsymbol{\theta}) = -\ln L(\boldsymbol{\theta}) + \lambda \Omega(\boldsymbol{\theta}) \tag{3}$$

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  $\lambda$  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 (3), représenté par la frontière de décision, pour la même régularisation  $\Omega(\boldsymbol{\theta})$  et différentes valeurs de  $\lambda$ . Discuter la qualité prédictive de chacun des trois modèles.

**Solution 3**

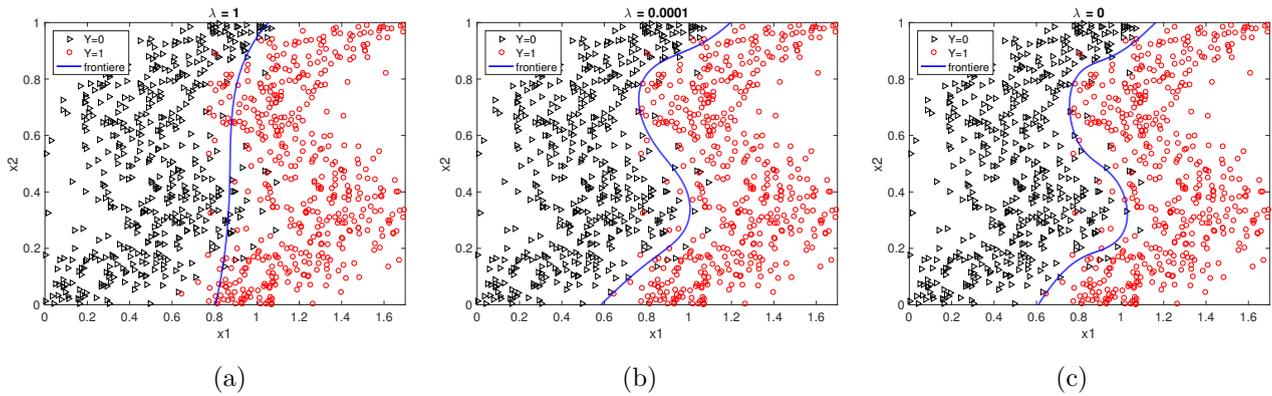


FIGURE 2 – Nuage de données ( $\circ, \blacktriangleright$ ) et modèle de prédiction (—)

1. Voir cours/td
2. Voir cours/td
3. On peut voir sur le graphique (a), pour lequel le niveau de pénalisation  $\lambda$  pour apprendre le modèle, est le plus fort parmi les trois situations, le modèle prédictif est un peu simpliste ; il s'agit d'un sous-apprentissage. Le modèle du graphique (c) correspond quant à lui à l'estimateur du maximum de vraisemblance non-régularisé ( $\lambda = 0$ ), et on peut voir qu'il sur-apprend un peu et devient spécifique aux données d'apprentissage. Le meilleur modèle est celui du graphique (b) et correspond au meilleur compromis entre l'attache aux données et la complexité du modèle. D'un point de vue décomposition biais-variance du risque quadratique, on a :
  - modèle (a) : grand biais mais petite variance,
  - modèle (c) : petit biais mais grande variance,
  - modèle (b) : biais et variance tous les deux petits.

**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 :

$$\hat{y}_i = \arg \max_{k \in \llbracket 0, 1 \rrbracket} \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \hat{\Psi}), \quad (4)$$

$\hat{\Psi}$  étant l'estimateur du vecteur paramètre  $\Psi$  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  $\mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = k; \boldsymbol{\theta}_k)$ .

2. On note par  $\alpha = \mathbb{P}(Y_i = 1)$ , la probabilité a priori de la classe 1. On a alors  $\boldsymbol{\Psi} = (\alpha, \boldsymbol{\theta}_1^\top, \boldsymbol{\theta}_0^\top)^\top$  que l'on cherche à estimer la méthode du maximum de vraisemblance. Donner l'expression de la fonction de vraisemblance

$$L(\boldsymbol{\Psi}) = \mathbb{P}((\mathbf{X}_1 = \mathbf{x}_1, Y_1 = y_1), \dots, (\mathbf{X}_n = \mathbf{x}_n, Y_n = y_n); \boldsymbol{\Psi}). \quad (5)$$

3. En déduire celle de la fonction de log-vraisemblance  $\ln L(\boldsymbol{\Psi})$ .
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  $\boldsymbol{\theta}$  par :

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

5. On considère le jeu de données de la Table 1. Prédire en appliquant la règle (4) la valeur de  $Y$  pour le dernier vecteur.

| $\mathbf{X}$ | $Y$ |   |
|--------------|-----|---|
| 0            | 1   | 1 |
| 1            | 0   | 1 |
| 1            | 0   | 0 |
| 0            | 1   | 0 |
| 0            | 1   | 1 |
| 1            | 1   | ? |

TABLE 1 – Jeu de données

#### Solution 4

1. Pour cela il suffit d'exploiter l'hypothèse  $\mathcal{H}$ , qui correspond à une hypothèse d'indépendance conditionnelle des variables. Par cette indépendance conditionnelle on peut donc écrire : étant le paramètre de la loi de Bernoulli de la variable binaire  $x_{ij}, j = 1 \dots, d$ . Cette loi conditionnelle est la loi de Bernoulli multivariée pour la classe  $k$

$$\begin{aligned} \mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = k; \boldsymbol{\theta}_k) &= \prod_{j=1}^p \mathbb{P}(X_{ij} = x_{ij} | Y_i = k; \theta_{kj}) \\ &= \prod_{j=1}^p \theta_{kj}^{x_{ij}} (1 - \theta_{kj})^{1-x_{ij}}. \end{aligned} \quad (6)$$

2. L'échantillon étant indépendant, la fonction de vraisemblance est donc donnée par :

$$\begin{aligned}
L(\Psi) &= \mathbb{P}((\mathbf{X}_1 = \mathbf{x}_1, Y_1 = y_1), \dots, (\mathbf{X}_n = \mathbf{x}_n, Y_n = y_n); \Psi) \\
&= \prod_{i=1}^n \mathbb{P}(\mathbf{X}_i = \mathbf{x}_i, Y_i = y_i; \Psi) \\
&= \prod_{i=1}^n \mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = y_i; \Psi) \mathbb{P}(Y_i = y_i).
\end{aligned} \tag{7}$$

Comme  $Y$  est binaire, on peut donc écrire

$$L(\Psi) = \prod_{i=1}^n [\mathbb{P}(Y_i = 1) \mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = 1; \theta_1)]^{Y_i} [\mathbb{P}(Y_i = 0) \mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = 0; \theta_0)]^{1-Y_i},$$

où  $\mathbb{P}(Y_i = 1) = 1 - \mathbb{P}(Y_i = 0) = \alpha$ . Ce qui donne, d'après la loi conditionnelle (6) :

$$L(\Psi) = \prod_{i=1}^n \left[ \alpha \prod_{j=1}^p \theta_{1j}^{x_{ij}} (1 - \theta_{1j})^{1-x_{ij}} \right]^{Y_i} \left[ (1 - \alpha) \prod_{j=1}^p \theta_{0j}^{x_{ij}} (1 - \theta_{0j})^{1-x_{ij}} \right]^{1-Y_i}. \tag{8}$$

3. La fonction de log-vraisemblance  $\ln L(\Psi)$  est donc donnée par

$$\begin{aligned}
\ln L(\Psi) &= \sum_{i=1}^n \ln \left\{ \left[ \alpha \prod_{j=1}^p \theta_{1j}^{x_{ij}} (1 - \theta_{1j})^{1-x_{ij}} \right]^{Y_i} \left[ (1 - \alpha) \prod_{j=1}^p \theta_{0j}^{x_{ij}} (1 - \theta_{0j})^{1-x_{ij}} \right]^{1-Y_i} \right\} \\
&= \sum_{i=1}^n \left\{ Y_i \ln \alpha + Y_i \sum_{j=1}^d [x_{ij} \ln \theta_{1j} + (1 - x_{ij}) \ln(1 - \theta_{1j})] \right. \\
&\quad \left. + (1 - Y_i) \ln(1 - \alpha) + (1 - Y_i) \sum_{j=1}^d [x_{ij} \ln \theta_{0j} + (1 - x_{ij}) \ln(1 - \theta_{0j})] \right\}.
\end{aligned} \tag{9}$$

4. La dérivée partielle de la log-vraisemblance par rapport à  $\alpha$  est donnée par

$$\frac{\partial \ln L(\Psi)}{\partial \alpha} = \sum_{i=1}^n \left( \frac{y_i}{\alpha} - \frac{1 - y_i}{1 - \alpha} \right) = \frac{1}{\alpha(1 - \alpha)} \sum_{i=1}^n (Y_i - \alpha).$$

En l'annulant par rapport à  $\alpha$  on trouve l'EMV de  $\alpha$  :

$$\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n Y_i.$$

Celle par rapport à  $\theta_{1j}$  est donnée par

$$\begin{aligned}
\frac{\partial \ln L(\Psi)}{\partial \theta_{1j}} &= \frac{\partial \left\{ \sum_{i=1}^n y_i \sum_{j=1}^d [x_{ij} \ln \theta_{1j} + (1 - x_{ij}) \ln(1 - \theta_{1j})] \right\}}{\partial \theta_{1j}} \\
&= \sum_{i=1}^n y_i \left( \frac{x_{ij}}{\theta_{1j}} - \frac{1 - x_{ij}}{1 - \theta_{1j}} \right) \\
&= \sum_{i=1}^n y_i \left( \frac{x_{ij}(1 - \theta_{1j})}{\theta_{1j}(1 - \theta_{1j})} - \frac{\theta_{1j}(1 - x_{ij})}{\theta_{1j}(1 - \theta_{1j})} \right) \\
&= \frac{\sum_{i=1}^n y_i x_{ij} - \sum_{i=1}^n y_i \theta_{1j}}{\theta_{1j}(1 - \theta_{1j})}.
\end{aligned}$$

En prenant la valeur  $\theta_{1j}$  qui annule cette dérivée, on trouve l'EMV de  $\theta_{1j}$  :

$$\hat{\theta}_{1j} = \frac{\sum_{i=1}^n Y_i X_{ij}}{\sum_{i=1}^n Y_i},$$

qui correspond à la moyenne des descripteurs  $X_{ij}$  de la classe 1, et l'expression vectorielle correspondante est

$$\hat{\boldsymbol{\theta}}_1 = \frac{1}{\sum_{i=1}^n Y_i} \sum_{i=1}^n Y_i \mathbf{X}_i,$$

ce qui correspond au vecteur moyen de tous les vecteurs  $\mathbf{x}_i$  de la classe 1.

De façon similaire, on trouve l'EMV de  $\boldsymbol{\theta}_0$  :

$$\hat{\boldsymbol{\theta}}_{0j} = \frac{1}{\sum_{i=1}^n (1 - Y_i)} \sum_{i=1}^n (1 - Y_i) X_{ij},$$

ce qui correspond à la moyenne de tous les descripteurs  $X_{ij}$  de la classe 0 et on a

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

qui correspond à la moyenne des individus de la classe 0.

5. La prédiction consiste à affecter l'individu  $\mathbf{x}_i$  à la classe  $\hat{y}_i$  maximisant la probabilité a posteriori :

$$\begin{aligned} \hat{y}_i &= \arg \max_{k \in [0,1]} \mathbb{P}(Y_i = k | \mathbf{X}_i = \mathbf{x}_i; \hat{\boldsymbol{\Psi}}) \\ &= \arg \max_{k \in [0,1]} \mathbb{P}(Y_i = k) \mathbb{P}(\mathbf{X}_i = \mathbf{x}_i | Y_i = k; \hat{\boldsymbol{\Psi}}) \\ &= \arg \max_{k \in [0,1]} \mathbb{P}(Y_i = k) \prod_{j=1}^p \hat{\theta}_{kj}^{x_{ij}} (1 - \hat{\theta}_{kj})^{1-x_{ij}} \end{aligned} \quad (10)$$

$\hat{\boldsymbol{\Psi}}$  étant l'EMV du vecteur paramètre  $\boldsymbol{\Psi}$ . Afin de pouvoir appliquer cette règle, nous devons calculer tout d'abord  $\hat{\boldsymbol{\Psi}} = (\hat{\alpha}, \hat{\boldsymbol{\theta}}_1^\top, \hat{\boldsymbol{\theta}}_0^\top)^\top$ . Selon les données de la Table 1, on trouve :

$$\begin{aligned} \hat{\alpha} &= \frac{1}{n} \sum_{i=1}^n Y_i = \frac{3}{5}; \\ \hat{\boldsymbol{\theta}}_{11} &= \frac{\sum_{i=1}^n Y_i X_{i1}}{\sum_{i=1}^n Y_i} = \frac{1}{3}; \quad \hat{\boldsymbol{\theta}}_{12} = \frac{\sum_{i=1}^n Y_i X_{i2}}{\sum_{i=1}^n Y_i} = \frac{2}{3}; \\ \hat{\boldsymbol{\theta}}_{01} &= \frac{\sum_{i=1}^n (1 - Y_i) X_{i1}}{\sum_{i=1}^n (1 - Y_i)} = \frac{1}{2}; \quad \hat{\boldsymbol{\theta}}_{02} = \frac{\sum_{i=1}^n (1 - Y_i) X_{i2}}{\sum_{i=1}^n (1 - Y_i)} = \frac{1}{2}. \end{aligned}$$

Avec  $\mathbf{x}_i = (1, 1)^\top$ , pour la classe  $Y_i = 1$ , on a

$$\hat{\alpha} \prod_{j=1}^p \hat{\theta}_{kj}^{x_{ij}} (1 - \hat{\theta}_{kj})^{1-x_{ij}} = \frac{3}{5} \left(\frac{1}{3}\right)^1 \left(1 - \frac{1}{3}\right)^0 \left(\frac{2}{3}\right)^1 \left(1 - \frac{2}{3}\right)^0 = \frac{2}{15}$$

et pour la classe  $Y_i = 0$  on a

$$(1 - \hat{\alpha}) \prod_{j=1}^p \hat{\theta}_{kj}^{x_{ij}} (1 - \hat{\theta}_{kj})^{1-x_{ij}} = \frac{2}{5} \left(\frac{1}{2}\right)^1 \left(1 - \frac{1}{2}\right)^0 \left(\frac{1}{2}\right)^1 \left(1 - \frac{1}{2}\right)^0 = \frac{1}{10}.$$

On peut donc en déduire que  $\hat{y}_i = 1$ .