

Examen d'analyse de données - Durée 1h30

Les documents de cours (transparents de cours, sujets de TP, TD, CTD et notes manuscrites) sont autorisés. Les trois exercices sont indépendants.

Exercice 1 : Classification bayésienne : 7.5 points

Soient deux classes C_1 et C_2 équiprobables dans \mathbb{R}^2 . Les observations de ces deux classes suivent une loi Gaussienne avec pour paramètres respectifs : \mathbf{m}_1 et \mathbf{m}_2 (vecteurs moyennes) et Σ_1 et Σ_2 (matrices de variance-covariance). On rappelle l'expression de la densité de probabilité conditionnelle multivariée pour la classe C_i (avec $i \in \{1, 2\}$), appelée vraisemblance, définie dans \mathbb{R}^d (avec ici $d = 2$) :

$$p(\mathbf{x}|C_i) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i) \right].$$

Questions

1. En posant $\mathbf{x} = [x_1, x_2]^T$, $\mathbf{m}_1 = [m_{11}, m_{12}]^T$, $\mathbf{m}_2 = [m_{21}, m_{22}]^T$, $\Sigma^{-1} = \begin{bmatrix} a & 0 \\ 0 & d \end{bmatrix}$ et en supposant que les matrices de covariances sont égales : $\Sigma_1 = \Sigma_2 = \Sigma$ et que les classes sont équiprobables, montrer que la frontière de décision entre les 2 classes C_1 et C_2 est donnée par l'équation suivante :

$$(\mathbf{m}_2 - \mathbf{m}_1)^T \Sigma^{-1} \mathbf{x} + C = 0$$

où C est une constante. Donner l'expression de C en fonction de \mathbf{m}_1 , \mathbf{m}_2 et Σ .

2. On considère un jeu de données dont les moyennes, matrices de covariance et probabilités sont définies par

$$\mathbf{m}_1 = [0, 2]^T, \Sigma_1 = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}, P(C_1) = P(C_2) = 0.5, \mathbf{m}_2 = [0, 0]^T, \Sigma_2 = \Sigma_1.$$

- Soient les deux points $\mathbf{p}_1 = [3, -2]^T$, $\mathbf{p}_2 = [3, 2]^T$ à classer. Pour chacun de ces points, calculer les valeurs des log-vraisemblances pour les deux classes et donner la classe attribuée.
- Calculer l'équation de la frontière de décision. Commenter.

Correction :

1. Comme les deux classes sont équiprobables et que leurs matrices de covariance sont identiques, la frontière de décision définie par $P(C_1|\mathbf{x}) = P(C_2|\mathbf{x})$ se réduit à (en prenant le logarithme des vraisemblances :

$$(\mathbf{x} - \mathbf{m}_1)^T \Sigma^{-1} (\mathbf{x} - \mathbf{m}_1) = (\mathbf{x} - \mathbf{m}_2)^T \Sigma^{-1} (\mathbf{x} - \mathbf{m}_2)$$

avec $\Sigma = \Sigma_1 = \Sigma_2$. En développant, on obtient

$$-2\mathbf{x}^T \Sigma^{-1} \mathbf{m}_1 + \mathbf{m}_1^T \Sigma^{-1} \mathbf{m}_1 = -2\mathbf{x}^T \Sigma^{-1} \mathbf{m}_2 + \mathbf{m}_2^T \Sigma^{-1} \mathbf{m}_2$$

soit

$$(\mathbf{m}_2 - \mathbf{m}_1)^T \Sigma^{-1} \mathbf{x} + C = 0$$

avec

$$C = \frac{1}{2} \mathbf{m}_1^T \Sigma^{-1} \mathbf{m}_1 - \frac{1}{2} \mathbf{m}_2^T \Sigma^{-1} \mathbf{m}_2.$$

2. Pour le point p_1 ,

$$p(p_1|C_1) = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} [3 \ -4] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ -4 \end{bmatrix} = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} \left(\frac{9}{4} + 16 \right).$$

De même pour la classe C_2 ,

$$p(p_1|C_2) = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} [3 \ -2] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ -2 \end{bmatrix} = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} \left(\frac{9}{4} + 4 \right)$$

Donc d'après la règle de Bayes avec équiprobabilité des classes, $p(p_1|C_1) < p(p_1|C_2)$ donc p_1 est classé dans la classe C_2 .

Pour le point p_2 ,

$$p(p_2|C_1) = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} [3 \ 0] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 0 \end{bmatrix} = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} \left(\frac{9}{4} \right)$$

De même pour la classe C_2 ,

$$p(p_2|C_2) = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} [3 \ 2] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \end{bmatrix} = -\log(2\pi) - \frac{1}{2} \log(4) - \frac{1}{2} \left(\frac{9}{4} + 4 \right)$$

Donc d'après la règle de Bayes avec équiprobabilité des classes, $p(p_2|C_2) < p(p_2|C_1)$ donc p_2 est classé dans la classe C_1 .

3. En reprenant la question 1, on obtient :

$$[0 \ -2] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \frac{1}{2} [0 \ 2] \begin{bmatrix} \frac{1}{4} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \end{bmatrix} = 0 \Leftrightarrow -2x_2 + \frac{1}{2} 4 = 0.$$

Donc l'équation de la frontière de décision est une droite horizontale d'équation $x_2 = 1$ (pas besoin de x_1 pour faire la classification). On remarquera qu'on pouvait trouver ce résultat sans calcul à l'aide de la distance aux barycentres. En effet la droite d'équation $x_2 = 1$ est la médiatrice du segment $[m_1, m_2]$.

Exercice 2 : Modélisation d'une réaction chimique par la méthode des moindres carrés : 5 points

Dans une réaction chimique, on souhaite modéliser l'évolution de la concentration d'un réactif en fonction du temps. On a mesuré expérimentalement :

Temps (s)	7	18	27	56
Concentration	32	28	25	18

Dans la suite, on notera $(T_i)_{1 \leq i \leq 4}$ la suite des temps considérés et $(C_i)_{1 \leq i \leq 4}$ la suite des concentrations mesurées.

Nous souhaitons effectuer une modélisation de la réaction par une réaction chimique à l'ordre 1, c'est-à-dire, si $C(t)$ désigne la concentration en fonction du temps, on a

$$\frac{dC(t)}{dt} = -\lambda C(t) \tag{1}$$

où λ représente la constante de réaction. Cette équation admet pour solution $C(t) = C_0 e^{-\lambda t}$ où C_0 représente la concentration initiale. On souhaite estimer les paramètres réels C_0 et λ .

Questions

1. Justifier que $\lambda > 0$.
2. Ecrivez matriciellement le problème aux moindres carrés linéaire à résoudre (MCO) permettant d'estimer les paramètres (C_0, λ) , c'est-à-dire définissez $\beta \in \mathbb{R}^2$, $A \in \mathbb{R}^{4 \times 2}$ et $\mathbf{b} \in \mathbb{R}^4$ tels que $\hat{\beta}_{OLS}$ soit la solution du problème suivant :

$$\min_{\beta \in \mathbb{R}^2} \|\mathbf{A}\beta - \mathbf{b}\|^2.$$

3. Donner la solution analytique de ce problème (on ne demande pas de calculer la solution numérique).

Correction

1. Il s'agit d'une équation différentielle du 1er ordre dont la solution est $C(t) = C_0 \exp(-\lambda t)$. La fonction exponentielle est strictement croissante. Or la concentration mesurée décroît strictement en fonction du temps. Donc il faut prendre $\lambda > 0$.
2. En prenant le logarithme de la solution obtenue à l'équation précédente, on obtient :

$$\forall i \in \{1, \dots, 4\}, \ln(C_i) = \ln(C_0) - \lambda T_i$$

que l'on peut écrire sous forme matricielle :

$$\begin{bmatrix} \ln(C_1) \\ \vdots \\ \ln(C_4) \end{bmatrix} = \begin{bmatrix} -T_1 & 1 \\ \vdots & \vdots \\ -T_4 & 1 \end{bmatrix} \begin{bmatrix} \lambda \\ \ln C_0 \end{bmatrix}.$$

On peut donc se ramener à un problème aux moindres carrés linéaires, en posant $\beta = \begin{bmatrix} \lambda \\ \ln(C_0) \end{bmatrix}$

$$\min_{\beta \in \mathbb{R}^2} \|\mathbf{A}\beta - \mathbf{b}\|^2$$

$$\text{avec } \mathbf{A} = \begin{bmatrix} -7 & 1 \\ -18 & 1 \\ -27 & 1 \\ -56 & 1 \end{bmatrix} \text{ et } \mathbf{b} = \begin{bmatrix} \ln(32) \\ \ln(28) \\ \ln(25) \\ \ln(18) \end{bmatrix}.$$

3. La solution des moindres carrés est définie par

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} = \mathbf{A}^+ \mathbf{b}.$$

Exercice 3 : Rugby ! - 7.5 points

On cherche à construire un arbre de décision permettant de décider si une équipe de rugby (par exemple, le Stade Toulousain) va gagner ou perdre le prochain match. Une base d'apprentissage a été construite en considérant les données suivantes qui récapitulent les conditions qui accompagnent les succès et les échecs de cette équipe de rugby.

Match à domicile	Ciel	Match précédent gagné ?	Match gagné ?
oui	Soleil	oui	oui
oui	Pluie	non	non
oui	Soleil	non	oui
non	Couvert	oui	oui
non	Pluie	oui	oui
non	Soleil	non	non

Questions

1. Déterminer l'indice de Gini associé à cette base d'apprentissage vis-à-vis des deux classes "Match gagné" et "Match perdu". **2 points**
2. Déterminer la variation de l'indice de Gini lorsqu'on coupe les données à l'aide des variables "Match à domicile", "Ciel" et "Match précédent gagné ?" (**1.5 point par variable**). En déduire la variable qui sera utilisée au premier niveau de l'arbre de décision. (**1 point**)

Correction

1. Indice de Gini de la base (ici c'est "Match gagné ?"), $i \in \{1, 2\}$ pour {oui, non}
 n =nbre d'occurrences totales (ici $n = 6$) et n_i = nbre d'occurrences "oui" ou "non".

$$\text{Gini}(Jouer) = \sum_{i=1}^2 \frac{n_i}{n} \left(1 - \frac{n_i}{n}\right) = 1 - \sum_{i=1}^2 \left(\frac{n_i}{n}\right)^2 = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = \frac{4}{9}.$$

2. (a) Indice de Gini de la variable "Ciel": 3 sous-ensembles

$\frac{n_{se}}{n} = p_i$ =proportion du sous-ensemble dans la variable

- i. sous-ensemble "Soleil" : $i \in \{1, 2\}$ pour {oui, non}, $n_{se_s} = 3$:

$$\text{Gini}(\text{Ciel} = \text{Soleil}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9}$$

- ii. sous-ensemble "Couvert" : $i \in \{1, 2\}$ pour {oui, non}, $n_{se_c} = 1$:

$$\text{Gini}(\text{Ciel} = \text{couvert}) = 1 - \left(\frac{1}{1}\right)^2 - \left(\frac{0}{1}\right)^2 = 0$$

- iii. sous-ensemble "Pluie" : $n_{se_p} = 2$

$$\text{Gini}(\text{Ciel} = \text{Pluie}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = \frac{1}{2}$$

On a alors

$$\text{Gini}(\text{Ciel}) = \frac{n_{se_s}}{n} \text{Gini}(\text{Soleil}) + \frac{n_{se_c}}{n} \text{Gini}(\text{Couvert}) + \frac{n_{se_p}}{n} \text{Gini}(\text{pluie}) = \left(\frac{3}{6} \times \frac{4}{9}\right) + \left(\frac{1}{6} \times 0\right) + \left(\frac{2}{6} \times \frac{1}{2}\right) = \frac{7}{18}$$

- (b) Indice de Gini de la variable "Match à domicile" : 2 sous-ensembles

- i. sous-ensemble "oui" : $i \in \{1, 2\}$ pour {oui, non}, $n_{se_o} = 3$:

$$\text{Gini}(\text{Match à domicile} = \text{oui}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9}.$$

- ii. sous-ensemble "non" : $n_{se_f} = 3$

$$\text{Gini}(\text{Match à domicile} = \text{non}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9}.$$

D'où

$$\text{Gini}(\text{Match à domicile}) = \frac{n_{se_o}}{n} \text{Gini}(\text{oui}) + \frac{n_{se_f}}{n} \text{Gini}(\text{non}) = \left(\frac{3}{6} \times \frac{4}{9}\right) + \left(\frac{3}{6} \times \frac{4}{9}\right) = \frac{4}{9}.$$

(c) Indice de Gini de “Match précédent gagné ?” : 2 sous-ensembles

i. sous ensemble “oui” : $i \in \{1, 2\}$ pour $\{\text{oui}, \text{non}\}$, $n_{se_f} = 3$:

$$\text{Gini}(\text{ Match précédent gagné} = \text{oui}) = 1 - \left(\frac{3}{3}\right)^2 - \left(\frac{0}{3}\right)^2 = 0.$$

ii. sous-ensemble “non” : $n_{se_F} = 3$

$$\text{Gini}(\text{Match précédent gagné} = \text{non}) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = \frac{4}{9}.$$

On en déduit

$$\text{Gini}(\text{Match précédent gagné}) = \frac{n_{se_f}}{n} \text{Gini}(\text{oui}) + \frac{n_{se_F}}{n} \text{Gini}(\text{non}) = \left(\frac{3}{6} \times 0\right) + \left(\frac{3}{6} \times \frac{4}{9}\right) = \frac{2}{9}.$$

Pour connaître la première variable utilisée au premier niveau de l’arbre CART, on maximise le gain défini par :

$$\text{Gain}(\text{Variable}) = \text{Gini}(\text{base}) - \text{Gini}(\text{variable})$$

(a) $\text{Gain}(\text{ciel}) = \frac{4}{9} - \frac{7}{18} = \frac{1}{18}$

(b) $\text{Gain}(\text{Match à domicile}) = \frac{4}{9} - \frac{4}{9} = 0$

(c) $\text{Gain}(\text{Match précédent gagné}) = \frac{4}{9} - \frac{2}{9} = \frac{2}{9}.$

Le gain est maximal pour la variable “Match précédent gagné” qui sera utilisée au premier niveau de l’arbre.