

# Classement de données binaires lorsque les populations d'apprentissage et de test sont différentes

Julien Jacques\*, Christophe Biernacki\*\*

\*Laboratoire de Statistiques et Analyse des Données,  
Université Pierre Mendès France,  
38040 Grenoble Cedex 9, France.  
julien.jacques@iut2.upmf-grenoble.fr,  
<http://www.julien.jacques2.free.fr>

\*\*Laboratoire Paul Painlevé UMR CNRS 8524,  
Université Lille I,  
59655 Villeneuve d'Ascq Cedex, France.  
christophe.biernacki@math.univ-lille1.fr

**Résumé.** L'analyse discriminante généralisée suppose que l'échantillon d'apprentissage et l'échantillon test, qui contient les individus à classer, sont issus d'une même population. Lorsque ces échantillons proviennent de populations pour lesquelles les paramètres des variables descriptives sont différents, l'analyse discriminante généralisée consiste à adapter la règle de classification issue de la population d'apprentissage à la population test, en estimant un lien entre ces deux populations. Ce papier étend les travaux existant dans un cadre gaussien au cas des variables binaires. Afin de relever le principal défi de ce travail, qui consiste à déterminer un lien entre deux populations binaires, nous supposons que les variables binaires sont issues de la discrétisation de variables gaussiennes latentes. Une méthode d'estimation puis des tests sur simulations sont présentés, et une application dans un contexte biologique illustre ce travail.

## 1 Introduction

L'analyse discriminante classique suppose que l'échantillon d'apprentissage et l'échantillon test, qui contient les individus à classer, sont issus d'une même population. Depuis les travaux de Fisher (1936), qui introduit une règle de discrimination linéaire entre deux groupes, de nombreuses évolutions ont été proposées (cf. McLachlan (1992) pour une revue). Toutes ces évolutions concernent la nature de la règle de discrimination : paramétrique, semi-paramétrique ou encore non paramétrique.

Une évolution alternative, introduite par Van Franeker et Ter Brack (1993) puis développée par Biernacki et al. (2002), considère le cas où l'échantillon d'apprentissage et l'échantillon test ne sont pas nécessairement issus d'une même population. Biernacki et al. définissent plusieurs modèles d'*analyse discriminante généralisée* dans un contexte gaussien, et les expérimentent sur une application biologique dans laquelle les deux populations sont des oiseaux de mer

d'une même espèce, mais d'origines géographiques différentes.

Mais dans beaucoup de domaines, comme les assurances ou la médecine, un grand nombre d'applications traite de données binaires. L'objectif de ce papier est d'étendre l'analyse discriminante généralisée, établie dans un contexte gaussien, au cas des données binaires. La différence entre les populations d'apprentissage et de test peut être géographique (comme dans l'application biologique précédemment citée), mais aussi temporelle ou tout autre.

La prochaine section présente les données et le modèle des classes latentes pour les deux populations d'apprentissage et de test. La section 3 fait l'hypothèse que les données binaires sont une discrétisation de variables continues latentes. Cette hypothèse permet d'établir un lien entre les deux populations, qui conduit à proposer huit modèles d'analyse discriminante généralisée pour données binaires. La section 4 traite de l'estimation des paramètres de ces modèles. Dans la section 5, des tests sur simulations puis une application dans un contexte biologique illustrent l'efficacité de l'analyse discriminante généralisée vis-à-vis de l'analyse discriminante classique et de la classification automatique. Finalement, la dernière section discute des possibles extensions de ces travaux.

## 2 Les données et le modèle des classes latentes

Les données consistent en deux échantillons : un étiqueté,  $S$ , issu d'une population  $P$ , et un non étiqueté,  $S^*$ , issu d'une population  $P^*$ . Les deux populations  $P$  et  $P^*$  peuvent être différentes.

L'échantillon d'apprentissage  $S$  est composé de  $n$  couples  $(x_1, z_1), \dots, (x_n, z_n)$ , réalisations indépendantes du couple aléatoire  $(\mathbf{X}, \mathbf{Z})$  de distribution :

$$X_{\mathbf{Z}^k=1}^j \sim \mathcal{B}(\alpha_{kj}) \quad \forall j = 1, \dots, d \quad \text{et} \quad \mathbf{Z} \sim \mathcal{M}(1, p_1, \dots, p_K). \quad (1)$$

En utilisant l'hypothèse d'indépendance conditionnelle des variables explicatives  $X^j$  ( $j = 1, \dots, d$ ) (Everitt (1984); Celeux et Govaert (1991)), la distribution de probabilité de  $\mathbf{X}$  s'écrit :

$$f(x^1, \dots, x^d) = \sum_{k=1}^K p_k \prod_{j=1}^d \alpha_{kj}^{x^j} (1 - \alpha_{kj})^{1-x^j}. \quad (2)$$

L'échantillon test  $S^*$  est quant à lui composé de  $n^*$  individus pour lesquels seules les variables explicatives  $x_1^*, \dots, x_{n^*}^*$  sont connues (les variables sont les mêmes que pour l'échantillon d'apprentissage). Les appartenances aux classes  $z_1^*, \dots, z_{n^*}^*$  sont inconnues. Nous considérons les couples  $(x_i^*, z_i^*)$  ( $i = 1, \dots, n^*$ ) comme des réalisations indépendantes du couple aléatoire  $(\mathbf{X}^*, \mathbf{Z}^*)$  de distribution analogue à  $(\mathbf{X}, \mathbf{Z})$  mais de paramètres différents, notés  $\alpha_{kj}^*$  et  $p_k^*$ .

L'objectif de la discrimination généralisée est alors d'estimer les  $n^*$  étiquettes  $z_1^*, \dots, z_{n^*}^*$  inconnues en utilisant l'information contenue à la fois dans  $S$  et dans  $S^*$ . Le défi consiste ainsi à trouver un lien entre les deux populations  $P$  et  $P^*$ .

### 3 Relation entre les populations d'apprentissage et de test

#### 3.1 Hypothèse gaussienne sous-jacente

Dans un contexte multi-normal, une relation stochastique linéaire entre  $P$  et  $P^*$  est non seulement justifiée (avec très peu d'hypothèses) mais aussi intuitive (Biernacki et al. (2002)). Dans le cas de données binaires, comme il ne semble pas exister de relation intuitive, une hypothèse supplémentaire est faite : nous supposons que les données binaires sont dues à la discrétisation de variables gaussiennes latentes. Cette hypothèse n'est pas nouvelle en statistique : nous pouvons citer comme exemple les travaux de Thurstone (1927), qui utilise cette hypothèse dans son modèle de jugement comparatif pour choisir entre deux stimulus, ou encore les travaux d'Everitt (1987), qui propose un algorithme de classification pour des données binaires, catégoriques et continues.

Nous supposons donc que les variables explicatives  $X_{|Z^k=1}^j$ , de distribution de Bernoulli  $\mathcal{B}(\alpha_{kj})$ , sont issues de la discrétisation de variables continues latentes  $Y_{|Z^k=1}^j$  conditionnellement indépendantes et de distribution normale  $\mathcal{N}(\mu_{kj}, \sigma_{kj}^2)$  :

$$X_{|Z^k=1}^j = \begin{cases} 0 & \text{si } \lambda_j Y_{|Z^k=1}^j < 0 \\ 1 & \text{si } \lambda_j Y_{|Z^k=1}^j \geq 0 \end{cases} \quad \text{pour } j = 1, \dots, d, \quad (3)$$

où  $\lambda_j \in \{-1, 1\}$  est introduit pour ne pas avoir à choisir quelle valeur de  $X^j$ , 0 ou 1, correspond à une valeur positive de  $Y^j$ . Le rôle de ce nouveau paramètre est d'éviter aux variables binaires d'hériter de l'ordre naturel induit par les variables continues.

Ainsi, nous pouvons déduire la relation suivante entre  $\alpha_{kj}$ , et  $\lambda_j$ ,  $\mu_{kj}$  et  $\sigma_{kj}$  :

$$\alpha_{kj} = p(X_{|Z^k=1}^j = 1) = \begin{cases} \Phi\left(\frac{\mu_{kj}}{\sigma_{kj}}\right) & \text{si } \lambda_j = 1 \\ 1 - \Phi\left(\frac{\mu_{kj}}{\sigma_{kj}}\right) & \text{si } \lambda_j = -1 \end{cases} \quad (4)$$

où  $\Phi$  est la fonction de répartition de la loi normale centrée réduite. Il est important de noter que l'hypothèse d'indépendance conditionnelle rend le calcul de  $\alpha_{kj}$  très simple. Sans cette hypothèse, il devient très complexe d'évaluer les intégrales multi-dimensionnelles induites par le calcul de  $\alpha_{kj}$ , surtout lorsque la dimension  $d$  du problème est grande, ce qui est souvent le cas avec des données binaires.

Comme pour la variable  $\mathbf{X}$ , nous supposons aussi que les variables binaires  $\mathbf{X}^*$  sont dues à la discrétisation de variables continues latentes  $\mathbf{Y}^*$ . Les équations sont les mêmes que (3) et (4), en changeant  $\alpha_{kj}$  en  $\alpha_{kj}^*$ ,  $\mu_{kj}$  en  $\mu_{kj}^*$  et  $\sigma_{kj}$  en  $\sigma_{kj}^*$ . Le paramètre  $\lambda_j^*$  est naturellement supposé égal à  $\lambda_j$ .

Dans un contexte gaussien, Biernacki et al. (2002) définissent la relation linéaire stochastique (5) entre les variables continues de la population  $P$  et celles de la population  $P^*$ , en émettant deux hypothèses relativement naturelles : la transformation entre  $P$  et  $P^*$  est  $\mathcal{C}^1$  et la  $j$ -ème composante  $Y_{|Z^k=1}^{*j}$  de  $\mathbf{Y}^*_{|Z^k=1}$  ne dépend que de la  $j$ -ème composante  $Y_{|Z^k=1}^j$  de  $\mathbf{Y}_{|Z^k=1}$ . Sous ces hypothèses, la seule transformation possible est la suivante :

$$\mathbf{Y}^*_{|Z^{*k}=1} \sim A_k \mathbf{Y}_{|Z^k=1} + \mathbf{b}_k, \quad (5)$$

où  $A_k$  est une matrice diagonale de  $\mathbb{R}^{d \times d}$  contenant les éléments  $a_{kj}$  ( $1 \leq j \leq d$ ) et  $\mathbf{b}_k$  est un vecteur de  $\mathbb{R}^d$ .

## Analyse discriminante généralisée sur données binaires

En appliquant cette relation à nos variables continues latentes, nous pouvons en déduire la relation suivante entre les paramètres  $\alpha_{kj}^*$  et  $\alpha_{kj}$  :

$$\alpha_{kj}^* = \Phi\left(\delta_{kj} \Phi^{-1}(\alpha_{kj}) + \lambda_j \gamma_{kj}\right), \quad (6)$$

où  $\delta_{kj} = \text{sgn}(a_{kj})$ ,  $\text{sgn}(t)$  désignant le signe de  $t$ , et où  $\gamma_{kj} = b_{kj}/(|a_{kj}| \sigma_{kj})$ .

Ainsi, conditionnellement au fait que les paramètres  $\alpha_{kj}$  sont connus (ils seront estimés en pratique), l'estimation des  $Kd$  paramètres continus  $\alpha_{kj}^*$  est obtenue par l'estimation des paramètres relatifs au lien entre  $P$  et  $P^*$  :  $\delta_{kj}$ ,  $\gamma_{kj}$  et  $\lambda_j$ . Le nombre de paramètres continus ( $\gamma_{kj}$ ) à estimer est ainsi  $Kd$ , ce qui est équivalent à estimer directement  $\alpha_{kj}^*$  sans utiliser la population  $P$ . Par conséquent, il est nécessaire de réduire ce nombre de paramètres continus à estimer.

Pour cela, nous introduisons un certain nombre de sous-modèles en imposant des contraintes sur la transformation entre les deux populations  $P$  et  $P^*$ .

### 3.2 Modèles de contraintes sur la relation

**Modèle  $M_1$**  :  $\sigma_{kj}$  et  $A_k$  sont libres et  $b_k = 0$  ( $k = 1, \dots, K$ ). Le modèle est :

$$\alpha_{kj}^* = \Phi\left(\delta_{kj} \Phi^{-1}(\alpha_{kj})\right) \quad \text{avec} \quad \delta_{kj} \in \{-1, 1\}.$$

Cette transformation correspond soit à l'identité, soit à une permutation des modalités de  $\mathbf{X}$ .

**Modèle  $M_2$**  :  $\sigma_{kj} = \sigma$ ,  $A_k = aI_d$  avec  $a > 0$ ,  $I_d$  la matrice identité de  $\mathbb{R}^{d \times d}$  et  $\mathbf{b}_k = \beta \mathbf{e}$ , avec  $\beta \in \mathbb{R}$  et  $\mathbf{e}$  le vecteur de dimension  $d$  composé seulement de 1 (la transformation est indépendante du groupe et de la dimension). Le modèle est :

$$\alpha_{kj}^* = \Phi\left(\Phi^{-1}(\alpha_{kj}) + \lambda'_j |\gamma|\right) \quad \text{avec} \quad \lambda'_j = \lambda_j \text{sgn}(\gamma) \in \{-1, 1\} \text{ et } |\gamma| \in \mathbb{R}^+.$$

L'hypothèse  $a > 0$  est faite pour avoir l'identifiabilité du modèle, et n'induit aucune restriction. La même hypothèse est faite pour les deux modèles suivants.

**Modèle  $M_3$**  :  $\sigma_{kj} = \sigma_k$ ,  $A_k = a_k I_d$ , avec  $a_k > 0$  pour tout  $1 \leq k \leq K$ , et  $\mathbf{b}_k = \beta_k \mathbf{e}$ , avec  $\beta_k \in \mathbb{R}$  (la transformation ne dépend que du groupe). Le modèle est :

$$\alpha_{kj}^* = \Phi\left(\Phi^{-1}(\alpha_{kj}) + \lambda'_{kj} |\gamma_k|\right) \quad \text{avec} \quad \lambda'_{kj} = \lambda_j \text{sgn}(\gamma_k) \in \{-1, 1\} \text{ et } |\gamma_k| \in \mathbb{R}^+.$$

**Modèle  $M_4$**  :  $\sigma_{kj} = \sigma_j$ ,  $A_k = A$ , avec  $a_{kj} > 0$  pour tout  $1 \leq k \leq K$  et  $1 \leq j \leq d$ , et  $\mathbf{b}_k = \beta$  avec  $\beta \in \mathbb{R}^d$  (la transformation ne dépend que de la dimension). Le modèle est :

$$\alpha_{kj}^* = \Phi\left(\Phi^{-1}(\alpha_{kj}) + \gamma'_j\right) \quad \text{avec} \quad \gamma'_j = \lambda_j \gamma_j \in \mathbb{R}.$$

Notons que dans ce modèle  $M_4$ , comme le paramètre  $\gamma'_j$  est libre, il inclut le paramètre  $\lambda_j$ .

Pour chacun des ces quatre modèles  $M_i$  ( $i = 1, \dots, 4$ ), nous prenons en compte une hypothèse supplémentaire sur les proportions des groupes ( $p_i$ ,  $i = 1, \dots, K$ ) : elles sont conservées

ou non de  $P$  vers  $P^*$ . Nous notons  $M_i$  le modèle avec proportions inchangées, et  $pM_i$  le modèle avec des proportions potentiellement changées. Huit modèles sont ainsi définis.

Notons que le modèle  $M_2$  est toujours inclus dans les modèles  $M_3$  et  $M_4$ , et  $M_1$  peut parfois être inclus dans les trois autres modèles.

Finalement, pour choisir automatiquement parmi ces huit modèles de discrimination généralisée, le critère BIC (*Bayesian Information Criterion*, Schwarz (1978)) peut être employé. Il est défini par :

$$\text{BIC} = -2l(\hat{\theta}) + \nu \log(n),$$

où  $\theta = \{p_k^*, \delta_{kj}, \lambda_j, \gamma_{kj}\}$  pour  $1 \leq k \leq K$  et  $1 \leq j \leq d$ ,  $l(\hat{\theta})$  est le maximum de la log-vraisemblance correspondant à l'estimation  $\hat{\theta}$  de  $\theta$ , et  $\nu$  est le nombre de paramètres continus libres associés au modèle donné (cf. tableau 1). Le modèle qui conduit à la plus petite valeur de BIC est retenu.

	$M_1$	$M_2$	$M_3$	$M_4$	$pM_1$	$pM_2$	$pM_3$	$pM_4$
continus	0	1	$K$	$d$	$K - 1$	$K$	$2K - 1$	$d + K - 1$
discrets	$Kd$	$d$	$Kd$	0	$Kd$	$d$	$Kd$	0

**TAB. 1** – Nombre de paramètres continus et discrets à estimer pour les huit modèles de discrimination généralisée.

Il est maintenant nécessaire d'estimer le paramètre  $\theta$ . La méthode du maximum de vraisemblance est retenue.

## 4 Estimation des paramètres de transition

L'analyse discriminante généralisée nécessite trois étapes d'estimation. Nous présentons la situation où les proportions sont inconnues au sein de  $P^*$ , le cas contraire étant immédiat.

La première étape consiste à estimer les paramètres  $p_k$  et  $\alpha_{kj}$  ( $1 \leq k \leq K$  et  $1 \leq j \leq d$ ) de la population  $P$  à partir de l'échantillon d'apprentissage  $S$ . Comme  $S$  est étiqueté, les estimateurs du maximum de vraisemblance sont connus et usuels (Everitt (1984); Celeux et Govaert (1991)).

La deuxième étape consiste à estimer les paramètres  $p_k^*$  et  $\alpha_{kj}^*$  ( $1 \leq k \leq K$  et  $1 \leq j \leq d$ ) du mélange de Bernoulli à partir de  $\theta$  et  $S^*$ . Pour estimer les  $\alpha_{kj}^*$ , nous devons estimer les paramètres du lien entre  $P$  et  $P^*$  : une fois  $\delta_{kj}$ ,  $\gamma_{kj}$  et  $\lambda_j$  estimés, un estimateur de  $\alpha_{kj}^*$  est déduit par l'équation (6). Cette étape est détaillée ci-dessous.

Enfin, la troisième étape consiste à estimer les appartenances aux groupes des individus de l'échantillon test  $S^*$ , et ce par maximum *a posteriori*.

Pour la deuxième étape, l'estimation par maximum de vraisemblance peut être réalisée à l'aide de l'algorithme EM (Dempster et al. (1977)) qui est bien adapté au cas des données manquantes (qui sont ici les appartenances aux classes).

## Analyse discriminante généralisée sur données binaires

La vraisemblance s'écrit :

$$L(\boldsymbol{\theta}) = \prod_{i=1}^{n^*} \sum_{k=1}^K p_k^* \prod_{j=1}^d \alpha_{kj}^{*x_i^{*j}} (1 - \alpha_{kj}^*)^{1-x_i^{*j}}.$$

La log-vraisemblance complétée est donnée par :

$$l_c(\boldsymbol{\theta}; z_1^*, \dots, z_{n^*}^*) = \sum_{i=1}^{n^*} \sum_{k=1}^K z_i^{*k} \log \left( p_k^* \prod_{j=1}^d \alpha_{kj}^{*x_i^{*j}} (1 - \alpha_{kj}^*)^{(1-x_i^{*j})} \right).$$

**L'étape E.** En utilisant une valeur courante  $\boldsymbol{\theta}^{(q)}$  du paramètre  $\boldsymbol{\theta}$ , l'étape E de l'algorithme EM consiste à calculer l'espérance conditionnelle de la log-vraisemblance complétée :

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)}) &= E_{\boldsymbol{\theta}^{(q)}} [l_c(\boldsymbol{\theta}; Z_1^*, \dots, Z_{n^*}^*) | x_1^*, \dots, x_{n^*}^*] \\ &= \sum_{i=1}^{n^*} \sum_{k=1}^K t_{ik}^{(q)} \left\{ \log(p_k^*) + \log \left( \prod_{j=1}^d \alpha_{kj}^{*x_i^{*j}} (1 - \alpha_{kj}^*)^{(1-x_i^{*j})} \right) \right\} \end{aligned}$$

où

$$t_{ik}^{(q)} = p(Z_i^{*k} = 1 | x_1^*, \dots, x_{n^*}^*; \boldsymbol{\theta}^{(q)}) = \frac{p_k^{*(q)} \prod_{j=1}^d (\alpha_{kj}^{*(q)})^{x_i^{*j}} (1 - \alpha_{kj}^{*(q)})^{(1-x_i^{*j})}}{\sum_{\kappa=1}^K p_{\kappa}^{*(q)} \prod_{j=1}^d (\alpha_{\kappa j}^{*(q)})^{x_i^{*j}} (1 - \alpha_{\kappa j}^{*(q)})^{(1-x_i^{*j})}}$$

est la probabilité conditionnelle que l'individu  $i$  appartienne au groupe  $k$ .

**L'étape M.** L'étape M de l'algorithme EM consiste à choisir la valeur de  $\boldsymbol{\theta}^{(q+1)}$  qui maximise l'espérance conditionnelle  $\mathcal{Q}$  calculée à l'étape E :

$$\boldsymbol{\theta}^{(q+1)} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)}). \quad (7)$$

Nous décrivons cette étape pour chaque composante de  $\boldsymbol{\theta} = \{p_k^*, \delta_{kj}, \lambda_j, \gamma_{kj}\}$ . Pour les proportions, cette maximisation donne l'estimateur suivant

$$p_k^{*(q+1)} = \frac{1}{n^*} \sum_{i=1}^{n^*} t_{ik}^{(q)}.$$

Pour les paramètres continus  $\gamma_{kj}$ , on montre pour chaque modèle que  $\mathcal{Q}$  est une fonction de  $\gamma_{kj}$  strictement concave et qui tend vers  $-\infty$  lorsqu'une norme du vecteur de paramètre  $(\gamma_{11}, \dots, \gamma_{Kd})$  tend vers  $\infty$  (cf. Jacques (2005)). Nous pouvons donc utiliser un algorithme de type Newton pour trouver l'unique maximum de  $\mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)})$  suivant  $\boldsymbol{\theta}$ .

Pour les paramètres discrets  $\delta_{kj}$  et  $\lambda_j$ , si la dimension  $d$  et le nombre de groupes  $K$  sont relativement petits (par exemple  $K = 2$  et  $d = 5$ ), la maximisation est faite en calculant  $\mathcal{Q}(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)})$  pour toutes les valeurs possibles des paramètres discrets. Si  $K$  ou  $d$  sont trop grands, le nombre de valeurs possibles de  $\delta_{kj}$  est trop important (par exemple  $2^{20}$  pour  $K = 2$  et  $d = 10$ ), et il

est donc impossible de parcourir toutes ces valeurs dans un temps raisonnable. Dans ce cas, nous utilisons une méthode de relaxation, qui consiste à supposer que le paramètre  $\delta_{kj}$  (respectivement  $\lambda_j$ ) n'est pas un paramètre binaire dans  $\{-1, 1\}$  mais continu dans  $[-1, 1]$ , noté  $\tilde{\delta}_{kj}$  (Wolsey (1998)). L'optimisation est alors faite sur le paramètre continu (avec un algorithme de type Newton comme pour  $\gamma_{kj}$  puisque  $\mathcal{Q}$  est une fonction de  $\tilde{\delta}_{kj}$  strictement concave), et la solution  $\tilde{\delta}_{kj}^{(q+1)}$  est discrétisée de la façon suivante pour obtenir une solution binaire  $\delta_{kj}^{(q+1)}$  :  $\delta_{kj}^{(q+1)} = \text{sgn}(\tilde{\delta}_{kj}^{(q+1)})$ .

**Remarque.** En pratique, le bouclage sur les valeurs possibles des paramètres discrets ne se fait pas au sein de l'étape M, mais en dehors de l'algorithme EM : nous estimons les paramètres  $p_k^*$  et  $\gamma_{kj}$  à l'aide de l'algorithme EM pour chaque valeur possible des paramètres discrets  $\delta_{kj}$  et  $\lambda_j$ , puis nous choisissons la solution de vraisemblance maximale.

## 5 Tests et application

### 5.1 Tests sur données simulées

Un grand nombre de tests sur simulations a été réalisé pour valider les huit modèles de discrimination généralisée, et nous en présentons un résumé dans cette section. L'objectif de ces tests est de comparer, suivant le taux de mauvais classement (taux d'erreur), l'analyse discriminante classique (DA), la classification automatique (aussi appelée *Clustering*, et qui revient à ne travailler que sur la population test en oubliant la population d'apprentissage), à l'analyse discriminante généralisée, pour laquelle nous avons sélectionné deux modèles, le meilleur modèle suivant le critère BIC (*GDA BIC*) et le meilleur modèle suivant le taux d'erreur (*GDA error*).

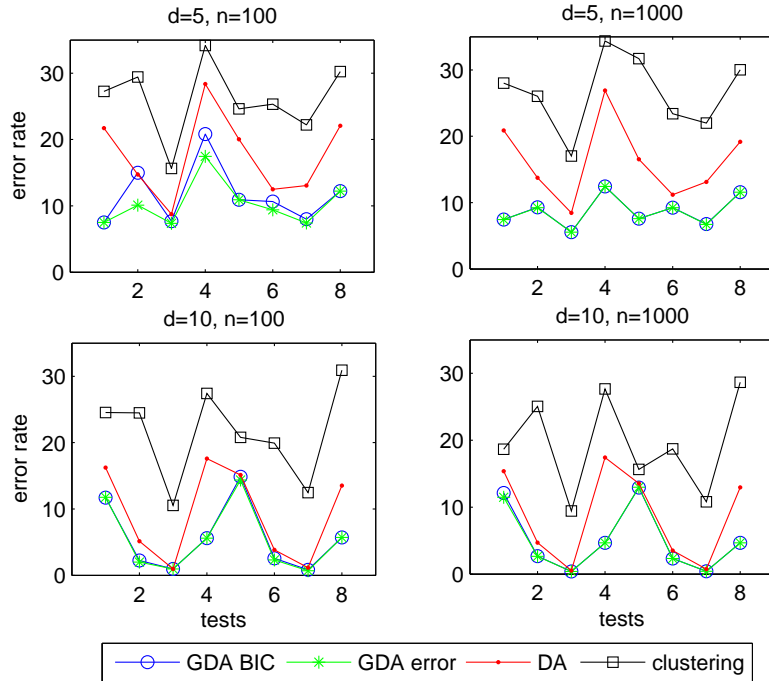
#### 5.1.1 Bon modèle

Dans une première série de tests, les données binaires sont simulées à partir de la discrétisation de données gaussiennes. De plus, les variables gaussiennes utilisées sont indépendantes conditionnellement à l'appartenance à un groupe, et la transformation entre  $P$  et  $P^*$  est choisie de sorte à être  $\mathcal{C}^1$  et de matrice d'homothétie  $A_k$  diagonale. Toutes les hypothèses de l'analyse discriminante généralisée sont ainsi respectées.

Les tests sont réalisés pour deux tailles d'échantillon ( $n = 100$  et  $n = 1000$ ), deux dimensions ( $d = 5$  et  $d = 10$ ), et pour huit transformations correspondant aux huit modèles de discrimination généralisée. Le détail de ces transformations est donné dans Jacques (2005).

Pour l'algorithme EM, la convergence est fixée à un gain en log-vraisemblance entre deux étapes de l'algorithme inférieur à  $10^{-6}$ , le nombre maximum d'itérations étant fixé à 200. L'initialisation des paramètres est faite de la façon suivante : 0 pour les paramètres  $\gamma$ ,  $\frac{1}{K}$  pour les proportions. Pour  $\lambda_j$  qui est égal à  $-1$  ou  $1$ , s'il est possible d'énumérer toutes ses valeurs possibles (nous rappelons que  $\lambda_j$  peut dépendre de la dimension et du nombre de groupes), nous choisissons la solution de log-vraisemblance maximale. Dans le cas contraire, nous utilisons la méthode de relaxation décrite précédemment et nous initialisons le paramètre  $\lambda_j$  à 0. De même pour le paramètre  $\delta_{kj}$  dans le modèle  $M_1$ .

## Analyse discriminante généralisée sur données binaires



**FIG. 1** – Tests sur données simulées sans bruit. La première ligne correspond à  $d = 5$  et la seconde à  $d = 10$ . La première colonne correspond à  $n = 100$  et la seconde à  $n = 1000$ .

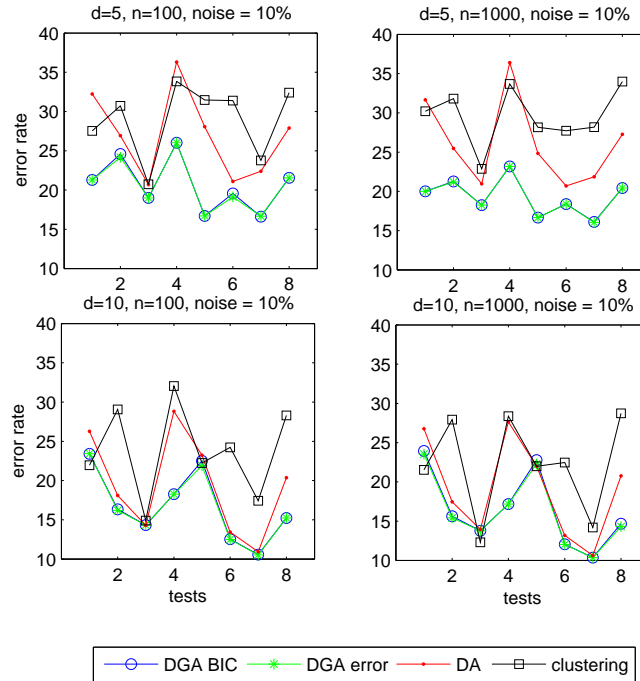
Les résultats sont présentés dans la figure 1. L'axe des abscisses représente les huit transformations entre  $P$  et  $P^*$  effectuées, les abscisses 1 à 4 correspondant au modèle  $pM_1$  à  $pM_4$  et les abscisses 5 à 8 correspondant au modèle  $M_1$  à  $M_4$  (cette convention sera conservée pour les prochaines figures 2 et 3). Ces résultats montrent que l'analyse discriminante généralisée donne des taux d'erreur plus faibles que l'analyse discriminante classique ou que la classification automatique. Ces bons résultats ne sont pas surprenant puisque toutes les hypothèses de l'analyse discriminante généralisée sont respectées.

Nous présentons maintenant quelques autres tests où ces hypothèses sont violées.

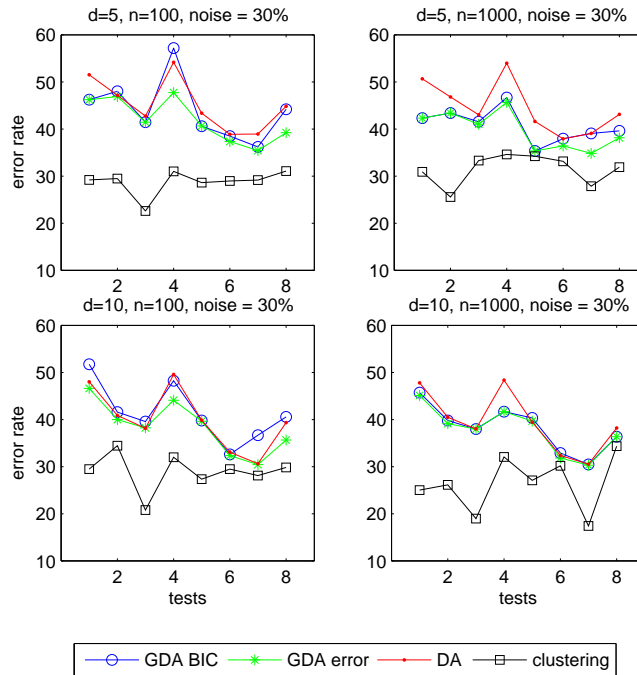
### 5.1.2 Mauvais modèles

Plusieurs autres séries de tests dans lesquelles les hypothèses de la discrimination généralisée ne sont pas vérifiées ont été réalisées. Nous présentons dans les figures 2 et 3 les résultats d'une de ces séries de tests, dans laquelle l'hypothèse que les variables binaires sont dues à la discrétisation de variables latentes gaussiennes est perturbée par l'introduction d'un bruit dans les données. Ce bruit consiste en une proportion de données uniformes parmi les données gaussiennes, centrées sur la moyenne des gaussiennes et étendue à plus ou moins deux écarts-types. Deux proportions de bruit sont testées : 10% et 30%.





**FIG. 2** – Tests sur données simulées avec 10% de bruit. La première colonne correspond à  $n = 100$  et la seconde à  $n = 1000$ . La première ligne correspond à  $d = 5$  et la seconde à  $d = 10$ .



**FIG. 3** – Tests sur données simulées avec 30% de bruit. La première colonne correspond à  $n = 100$  et la seconde à  $n = 1000$ . La première ligne correspond à  $d = 5$  et la seconde à  $d = 10$ .

## Analyse discriminante généralisée sur données binaires

Pour 10% de bruit (figure 2), l'analyse discriminante généralisée est encore meilleure que les autres méthodes, mais lorsque le bruit est trop important (figure 3), la classification automatique devient la meilleure méthode (suivant le taux de mauvais classement).

Trois autres séries de tests violant les hypothèses de la discrimination généralisée ont également été réalisées. Dans la première série, les données gaussiennes, à partir desquelles sont simulées les données binaires par discrétisation, sont bruitées par l'introduction de données binaires uniformément distribuées. Dans la deuxième série, les données binaires sont simulées par discrétisation de variables gaussiennes non conditionnellement indépendantes. Finalement, dans une troisième série, la matrice d'homothétie  $A_k$  de la transformation entre populations est choisie non diagonale. Dans tous ces tests, les résultats sont similaires à ceux présentés précédemment, c'est-à-dire lorsque l'écart au modèle n'est pas trop important, l'analyse discriminante généralisée est meilleure que les autres méthodes, mais lorsque cet écart est trop important, c'est la classification automatique qui doit être préférée.

Tous ces tests nous permettent de conclure que l'analyse discriminante généralisée donne de meilleurs classements que l'analyse discriminante classique ou que la classification automatique lorsque les populations d'apprentissage et de test ne sont pas identiques, et qu'elle est relativement robuste aux hypothèses suivantes : les données binaires sont issues de la discrétisation de variables latentes gaussiennes, et la transformation linéaire entre  $P$  et  $P^*$  s'effectue avec une matrice d'homothétie  $A_k$  diagonale.

Nous présentons maintenant une application sur un jeu de données réelles dans un contexte biologique.

### 5.2 Application sur données réelles

Les premières motivations pour lesquelles l'analyse discriminante généralisée a été développée sont des applications biologiques (Biernacki et al. (2002); Van Franeker et Ter Brack (1993)), dans lesquelles l'objectif est de prédire le sexe d'oiseaux à partir de variables biométriques. De très bons résultats ont été obtenus sous une hypothèse multi-normale.

L'espèce d'oiseaux considérée est l'espèce *Calanectris diomedea* (Thibault et al. (1997)). Cette espèce peut être divisée en trois sous-espèces, parmi lesquelles les *borealis*, qui vivent dans des îles de l'Atlantique (Açores, Canaries, etc.), et les *diomedea*, qui vivent dans des îles de la Méditerranée (Baléares, Corse, etc.). Les oiseaux de la sous-espèce *borealis* sont généralement plus grands que ceux de la sous-espèce *diomedea*. C'est pourquoi l'analyse discriminante généralisée n'est pas adaptée pour prédire le sexe des oiseaux *diomedea* à partir d'un échantillon d'apprentissage issu de la population des *borealis*.

La figure 4 illustre les différences entre deux variables morphologiques (tailles des ailes et de la queue), pour les deux espèces *diomedea* et *borealis*.

Un échantillon de *borealis* ( $n = 206$ , 45% de femelles) est constitué à partir d'information fournies par plusieurs musées nationaux. Cinq variables morphologiques sont mesurées : profondeur et longueur du bec, longueur du tarse (os de la patte), longueur des ailes, longueur de la queue. De même, un échantillon de *diomedea* ( $n = 38$ , 58% de femelles) est mesuré sur le même ensemble de variables morphologiques. Dans cet exemple deux groupes sont présents,

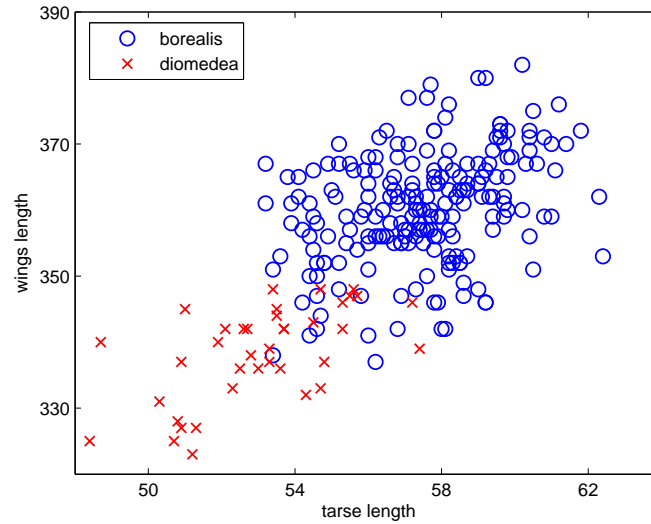


FIG. 4 – Longueur des ailes et de la queue pour les deux espèces *diomedea* et *borealis*.

les mâles et les femelles, et tous les oiseaux sont de sexe connus (par dissection). Pour procurer une application à notre travail, nous discrétisons les variables biométriques en variables binaires (petites ou grandes ailes, ...). Nous choisissons pour population d'apprentissage l'espèce *borealis* et pour population test l'espèce *diomedea*. Les trois méthodes de classification, discrimination généralisée, discrimination classique et classification automatique sont testées. Les résultats sont présentés dans le tableau 2.

	$pM_1$	$pM_2$	$pM_3$	$pM_4$	$M_1$	$M_2$	$M_3$	$M_4$	DA	Clustering
Taux	57.9	26.3	23.7	21.0	57.9	23.7	<b>15.8</b>	18.4	42.1	23.7
Nombre	22	10	9	8	22	9	<b>6</b>	7	16	9
BIC	269.7	222.5	<b>220.5</b>	237.0	267.3	221.6	221.5	233.6	648.4	-

TAB. 2 – Taux (en %) et nombre de mauvais classements, et valeurs du critère BIC pour la population test *diomedea* avec la population d'apprentissage *borealis*.

Si nous comparons les résultats obtenus en terme de taux d'erreur, la discrimination généralisée est la meilleure méthode avec un taux d'erreur de 15.8% pour le modèle  $M_3$ . Cette erreur est plus faible que celles obtenues par discrimination classique (42.11%) et par classification automatique (23.7%). Si l'on utilise le critère BIC pour choisir un modèle, trois sont mis en valeur ( $pM_3$ ,  $M_3$  et  $M_2$ ) parmi lesquels celui qui conduit au taux d'erreur le plus faible ( $M_3$ ). Cette application illustre l'intérêt de la discrimination généralisée vis-à-vis des procédures classiques de discrimination linéaire ou de classification automatique. En effet, en adaptant la

règle de classement issue de l'échantillon d'apprentissage à l'échantillon test en fonction des différences entre les deux populations d'apprentissage et de test, la discrimination généralisée permet de classer les individus de la population test de façon plus efficace qu'en appliquant directement la règle de classement issue de l'échantillon d'apprentissage, ou encore en oubliant l'échantillon d'apprentissage et en effectuant une classification automatique directement sur l'échantillon test.

Remarquons aussi que l'hypothèse gaussienne conditionnellement au sexe était relativement acceptable sur les variables biométriques sous-jacentes. Néanmoins, il y avait une forte corrélation entre ces différents caractères, violant alors les hypothèses d'indépendance conditionnelle de nos modèles.

## 6 Conclusion

L'analyse discriminante généralisée étend l'analyse discriminante classique en permettant aux échantillons d'apprentissage et de test d'être issus de populations partiellement différentes. Notre contribution consiste à adapter les travaux précurseurs réalisés dans un contexte gaussien au cas des données binaires. Une application en biologie illustre ce travail. En utilisant les modèles de discrimination généralisée définis dans ce papier, nous apportons un classement des oiseaux suivant leur sexe meilleur que ceux obtenus par analyse discriminante classique ou par classification automatique. L'application de ces méthodes sur des données issues de compagnie d'assurance est notre prochain objectif.

Les perspectives méthodologiques de ces travaux sont nombreuses.

Tout d'abord, nous avons défini le lien entre les deux populations en utilisant la fonction de répartition de la loi normale centrée réduite. Bien qu'il put paraître initialement difficile de trouver un lien entre populations binaires, un lien simple a été obtenu. Il n'aurait pas été facile de l'imaginer, mais il est aisément compréhensible *a posteriori*. Cela pourrait être intéressant d'essayer d'autres types de fonction de répartition (les raisons théoriques devront être développées et des tests devront être effectués).

Grâce à cette contribution, l'analyse discriminante généralisée est maintenant développée pour des données continues et des données binaires. Pour permettre de traiter un maximum de cas pratiques, il est important d'étudier le cas de variables catégoriques (à plus de deux modalités), puis ensuite le cas de variables mélangées (variables binaires, continues et catégoriques dans un même problème).

Finalement, il sera aussi très intéressant d'étendre les autres méthodes de discrimination classique comme la discrimination non-paramétrique.

## Références

- Biernacki, C., F. Beninel, et V. Bretagnolle (2002). A generalized discriminant rule when training population and test population differ on their descriptive parameters. *Biometrics* 58,2, 387–397.
- Celeux, G. et G. Govaert (1991). Clustering criteria for discrete data and latent class models. *Journal of classification* 8, 157–176.
- Dempster, A., N. Laird, et D. Rubin (1977). Maximum likelihood from incomplete data (with discussion). *Journal of the Royal Statistical Society Series B* 39, 1–38.
- Everitt, B. (1984). *An introduction to latent variables models*. London: Chapman & Hall.
- Everitt, B. (1987). A finite mixture model for the clustering of mixed-mode data. *Statistics and Probability Letters* 6, 305–309.
- Fisher, R. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics* 7, 179–188.
- Jacques, J. (2005). *Contributions à l'analyse de sensibilité et à l'analyse discriminante généralisée*. Thèse de doctorat, Université Joseph Fourier de Grenoble.
- McLachlan, G. (1992). *Discriminant Analysis and Statistical Pattern Recognition*. New-York: Wiley.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics* 6, 461–464.
- Thibault, J.-C., V. Bretagnolle, et C. Rabouam (1997). Cory's shearwater calonectris diomedea. *Birds of Western Palearctic Update 1*, 75–98.
- Thurstone, L. (1927). A law of comparative judgement. *Amer. J. Psychol.* 38, 368–389.
- Van Franeker, J. et C. Ter Brack (1993). A generalized discriminant for sexing fulmarine petrels from external measurements. *The Auk* 110(3), 492–502.
- Wolsey, L. (1998). *Integer Programming*. Wiley.

## Summary

Standard discriminant analysis supposes that both the training labelled sample and the test sample which has to be classed are issued from the same population. When these samples are issued from populations for which descriptive parameters are different, generalized discriminant analysis allows us to adapt the classification rule issued from the training population to the test population, by estimating a link between this two populations. This paper extends existing work available in a multi-normal context to the case of binary data. To raise the major challenge of this work which is to define a link between the two binary populations, we suppose that binary data are issued from the discretization of latent Gaussian data. Estimation method and tests on simulated data are presented, and an application in biological context illustrates this work.