# A DISCRIMINANT RULE UNDER TRANSFORMATION

Santiago Velilla and Juan A. Barrio[*]

Abstract ───────────────────────────────

We present a new rule for discriminating among continuous populations which are not multivariate normal. The basic idea is to construct the sample maximum likelihood discriminant rule after transforming the data by a suitable multivariate transformation to normality

Key words:
Cross-validation; Multivariate Box-Cox Transformation.

[*]Velilla, Departamento de Estadística y Econometría, Universidad Carlos III de Madrid; Barrio, Departamento de Física Atómica, Molecular y Nuclear, Universidad Complutense de Madrid.

# 1. INTRODUCTION

Discriminant analysis is a widely applied technique in multivariate data analysis. Given $g$ groups, populations or classes, $G_1$, $G_2$, ..., $G_g$, with corresponding probability density functions $f_i(x)$ on $\mathbb{R}^p$, $i=1$, ...,$g$, the problem is allocating an individual to one of these groups on the basis of his measurements $x=(x_1, ...,x_p)'$ which are considered as an observation from a population described by a random vector $X=(X_1,X_2, ...,X_p)'$ such that $X \sim f_i$ under $G_i$, $i=1$, ...,$g$. The allocation should be "optimal" in the sense of minimizing, on average, the number of incorrect assignments. When the densities $f_i(x)$ are known, a suitable discriminating rule is the *maximum likelihood* (ML) *discriminant rule* (Mardia et al. (1979), p. 301), defined as follows: Assign $x$ to $G_i$ if

$$f_i(x)= \max_{1 \le j \le g} f_j(x). \tag{1}$$

We will assume, for the rest of this work, that the a priori probabilities of the classes are equal. When this is the case, it can be shown that the rule (1) minimizes the total probability of misclassification (see Seber (1984), p. 331). When the densities $f_i(x)$ depend on unknown parameters $\theta_i$, i.e., if $f_i(x)=f_{\theta_i}(x)$, $i=1$, ...,$g$, we use the *sample ML* (SML) *discriminant rule* (Mardia et al. (1979), p. 309), which is obtained by replacing, in rule (1), the parameters $\theta_i$ by efficient estimates $\hat{\theta}_i=\hat{\theta}_i(\mathcal{X}_i)$, where $\mathcal{X}_i$ is a $n_i \times p$ data matrix from $G_i$, $i=1$, ...,$g$.

An important example of the ideas above appears when $f_i \sim N_p(\mu_i,\Sigma_i)$, and the parameters $(\mu_i,\Sigma_i)$ are unknown, $i=1$, ...,$g$. Let $\bar{x}_i$ and $S_i$ be, respectively, the sample mean and the covariance matrix of $\mathcal{X}_i$, $i=1$, ...,$g$. The SML rule can adopt two possible different forms: (*i*) If $\Sigma_1=\Sigma_2= ...=\Sigma_g=\Sigma$, say, we assign $x$ to $G_i$ if

1

$$\ell_i'x + c_i = \max_{1 \le j \le g} [\ell_j'x + c_j], \tag{2}$$

where $\ell_i = S_p^{-1}\overline{x}_i$, $c_i = -(1/2)\overline{x}_i'S_p^{-1}\overline{x}_i$, $i=1, \ldots, g$, and $S_p$ is the pooled estimator $S_p = \sum_{i=1}^{g} (n_i - 1)S_i/(N-g)$, $N = \sum_{i=1}^{g} n_i$; (ii) If the dispersion matrices are unequal, we assign $x$ to $G_i$ if

$$Q_i(x) = \min_{1 \le j \le g} Q_j(x), \tag{3}$$

where $Q_i(x) = (1/2)\log[|S_i|] + (1/2)[(x - \overline{x}_i)'S_i^{-1}(x - \overline{x}_i)]$, $i=1, \ldots, g$.

It is well-known that, when $g=2$, the rule (2) is equivalent to the *Fisher's Linear Discriminant Function* obtained by Fisher (1936) not assuming normality but looking for a "sensible" rule for discriminating between the groups. For $g \ge 2$, if we assume that the possible distributions of the random vector $X = (X_1, X_2, \ldots, X_p)'$ are continuous but not multivariate normal, the rules (2) and (3) above are, in general, no longer appropriate, and alternative procedures should be used. For example (see, e.g., chapter 6 in Seber (1984) or chapter 13 in Krzanowski (1990) for details) we could use a logistic discriminant rule or we could use the ML discriminant rule (1) replacing $f_i(x)$ by a nonparametric density estimator computed from $\mathcal{X}_i$, $i=1, \ldots, g$.

In this paper, we present a new method for discriminating among continuous populations which are not multivariate normal. The method is based on the multivariate generalization of the transformation of Box and Cox (1964) given by Andrews et al. (1971). In section 2, we introduce the general ideas of this new discrimination procedure. In section 3, we explicitly construct the associated discriminating rule and suggest a cross-validation method for assessing its performance. Section 4 contains an example of application and section 5 some final comments.

2

# 2. MOTIVATION

Let X be a random variable which takes values denoted by x. If X is not normal, it is often convenient to consider a transformation that might help to normalize the data. A useful family of transformations is the family of Box and Cox (1964):

$$x^{(\lambda)} = \begin{cases} \dfrac{x^{\lambda}-1}{\lambda} , & \lambda \neq 0 \\ \log x , & \lambda = 0. \end{cases} \qquad (4.1)$$

We assume that the following model holds

$$x^{(\lambda)} \sim N(\mu, \sigma^2), \qquad (4.2)$$

at least approximately. When we have a random vector $X=(X_1, X_2, \ldots, X_p)'$ whose distribution is not $N_p(\mu_1, \Sigma_1)$, Andrews et al. (1971) have given the following multivariate generalization of the setup (4.1)-(4.2). We have a p-vector $\Lambda = (\lambda_1, \ldots, \lambda_p)'$ of transformation parameters, one for each dimension, such that, approximately, the model

$$X^{(\Lambda)} = (X_1^{(\lambda_1)}, X_2^{(\lambda_2)}, \ldots, X_p^{(\lambda_p)})' \sim N_p(\mu, \Sigma) \qquad (5)$$

holds. In both (4) and (5), it is assumed that x and the components $X_j$ are positive. If not, x and $X_j$ must be shifted by adding suitable constants.

Write the parameters in model (5) as $\Theta = (\Lambda, \mu, \Sigma)$. Given an observation $x = (x_1, \ldots, x_p)'$ from the model (5), the likelihood of x is given by

$$\ell(x;\Theta) = (2\pi)^{-(p/2)} |\Sigma|^{-(1/2)} \exp[-(1/2)(x^{(\Lambda)}-\mu)'\Sigma^{-1}(x^{(\Lambda)}-\mu)] J_\Lambda(x), \qquad (6)$$

where $J_\Lambda(x)$ is the jacobian of the transformation, $J_\Lambda(x) = \prod_{j=1}^{p} x_j^{\lambda_j - 1}$. If we now face the problem of discriminating among groups $G_1, \ldots, G_g$ such that the respective densities $f_i(x)$ are not multivariate normal but can be, in turn, modelled by (5), i.e. such that $f_i(x) = f_{\Theta_i}(x) = \ell(x;\Theta_i)$, where $\Theta_i = (\Lambda_i, \mu_i, \Sigma_i)$, $i=1, \ldots, g$, we can construct the following SML

3

discriminant rule: Assign $\mathbf{x}$ to $G_i$ if

$$\ell(\mathbf{x};\hat{\theta}_i) = \max_{1 \le j \le g} \ell(\mathbf{x};\hat{\theta}_j), \tag{7}$$

where $\hat{\theta}_i = \hat{\theta}_i(\mathcal{X}_i)$ is an efficient estimator of $\Theta_i$, $i=1, \ldots, g$. In the next section, we show how to compute explicitly the discriminant rule (7), we relate it to the rules (2) and (3), and suggest a procedure for assessing its performance.

## 3. CONSTRUCTION AND ASSESSMENT OF THE DISCRIMINANT RULE

### 3.1 Two groups

We will consider first the case when we have only two groups, $G_1$ and $G_2$. For $i=1,2$, we are given a data matrix $\mathcal{X}_i$ of $n_i \times p$ from $G_i$ such that for some unknown vector $\Lambda_i = (\lambda_{i1}, \ldots, \lambda_{ip})'$ in $\mathbb{R}^p$ the rows of the $n_i \times p$ transformed data matrix $\mathcal{X}_i^{(\Lambda_i)} = (x_{ijk}^{(\lambda_{ik})})$, $j=1, \ldots, n_i$, $k=1, \ldots, p$, are i.i.d. $N_p(\mu_i, \Sigma_i)$. We also assume independence between $\mathcal{X}_1$ and $\mathcal{X}_2$.

### 3.1.1 Parameter estimation

Standard normal likelihood theory shows (see, e.g. Gnanadesikan (1977), p. 141), that, for $i=1,2$, the maximum likelihood estimator (MLE) of the parameter $\Theta_i = (\Lambda_i, \mu_i, \Sigma_i)$, denoted by $\hat{\theta}_i = (\hat{\Lambda}_i, \hat{\mu}_i, \hat{\Sigma}_i)$, is given by $\hat{\mu}_i = \bar{x}_i^{(\hat{\Lambda}_i)}$, $\hat{\Sigma}_i = ((n_i-1)/n_i) S_i^{(\hat{\Lambda}_i)}$, where $\bar{x}_i^{(\hat{\Lambda}_i)}$ and $S_i^{(\hat{\Lambda}_i)}$ are, respectively, the sample mean and the covariance matrix for $\mathcal{X}_i^{(\hat{\Lambda}_i)}$. For $i=1,2$, $\hat{\Lambda}_i$ is obtained by maximizing the corresponding concentrated log-likelihood which is (up to an additive constant)

$$L_{max,i}(\Lambda_i) = -(n_i/2)\log[\,|\mathcal{X}_i^{(\Lambda_i)'} Q_i \mathcal{X}_i^{(\Lambda_i)}|\,] + \log[J_{\Lambda_i}(\mathcal{X}_i)], \tag{8}$$

where $\qquad Q_i = (I_{n_i} - 1_{n_i} 1_{n_i}'/n_i) \qquad$ and $\qquad J_{\Lambda_i}(\mathcal{X}_i) = \prod_{k=1}^{p} J_{\lambda_{ik}},$

4

$J_{\lambda_{ik}} = \prod_{j=1}^{n_i} |d\alpha_{ijk}^{(\lambda_{ik})}/d\alpha_{ijk}| = (\prod_{j=1}^{n_i} \alpha_{ijk})^{\lambda_{ik}-1}$, k=1, ...,p, are the jacobian terms. Define, for each column k of $\mathcal{X}_i$, the $n_i$ vector of normalized variables $Z_{ik}^{(\lambda_{ik})} = (z_{ijk}^{(\lambda_{ik})})$, k=1, ...,p, of jth coordinate $z_{ijk}^{(\lambda_{ik})} = \alpha_{ijk}^{(\lambda_{ik})}/J_{\lambda_{ik}}^{1/n_i}$, j=1, ...,$n_i$. We have the following result.

**LEMMA 3.1** For i=1,2, let $Z_i^{(\Lambda_i)} = (Z_{i1}^{(\lambda_{i1})}, ...,Z_{ip}^{(\lambda_{ip})})$ be the $n_i \times p$ matrix of normalized variables associated with $\mathcal{X}_i$. We can write

$$L_{max,i}(\Lambda_i) = -(n_i/2)\log[\,|Z_i^{(\Lambda_i)'} Q_i Z_i^{(\Lambda_i)}|\,]. \tag{9}$$

**Proof.** Define, for i=1,2, the matrix $D_i = diag(J_{\lambda_{ik}}^{1/n_i}, ...., J_{\lambda_{ip}}^{1/n_i})$. (9) follows from (8), and by observing that $-(n_i/2)\log[J_{\Lambda_i}^{-2/n_i}(\mathcal{X}_i)] = -(n_i/2)\log[\,|D_i|^{-2}]$, and that $Z_i^{(\Lambda_i)} = \mathcal{X}_i^{(\Lambda_i)} D_i^{-1}$. ∎

For computational purposes, expression (9) is more convenient than (8) in order to determine the MLE $\hat{\Lambda}_i$. In experience of the authors, the functions $-L_{max,i}(\Lambda_i)$ are typically convex, a convenient feature for solving the optimization problem by means of a canned numerical routine.

### 3.1.2 Expression of the rule

Given a new individual $x=(x_1, ...,x_p)'$ define the p vector $u(x)=-(\log x_1, ..., \log x_p)'$. Two possible different situations arise:

(i) If $\Sigma_1 \neq \Sigma_2$, our proposal is estimating $\theta_i=(\Lambda_i,\mu_i,\Sigma_i)$ with $(\hat{\Lambda}_i,\hat{\mu}_i,\hat{S}_i)$, where $\hat{S}_i = S_i^{(\hat{\Lambda}_i)}$. The rule (7) can be written equivalently in the form: Assign x to $G_1$ if

$$\hat{Q}_1(x)+\hat{\Lambda}_1'u(x) = \min_{1 \leq j \leq 2} [\hat{Q}_j(x)+\hat{\Lambda}_j'u(x)], \tag{10}$$

5

where $\hat{Q}_i(x)=(1/2)\log[|\hat{S}_i|]+(1/2)[(x^{(\hat{\Lambda}_i)}-\hat{\mu}_i)'(\hat{S}_i)^{-1}(x^{(\hat{\Lambda}_i)}-\hat{\mu}_i)]$, i=1,2. We can interpret (10) as the estimated normal theory discriminant rule (3) applied to the transformed data matrices $\mathcal{X}_1^{(\hat{\Lambda}_1)}$ and $\mathcal{X}_2^{(\hat{\Lambda}_2)}$ and shifted by a linear term in u(x) which reflects the effect of the transformation;

(ii) If $\Sigma_1=\Sigma_2=\Sigma$, it is reasonable to estimate $\Sigma$ by the pooled estimator $\hat{S}_p = [(n_1-1)\hat{S}_1+(n_2-1)\hat{S}_2]/(N-2)$, where $N=n_1+n_2$. The rule (7) becomes now: Assign x to $G_1$ if

$$\tilde{Q}_i(x)+\hat{\Lambda}_i'u(x)= \min_{1 \leq j \leq 2} [\tilde{Q}_j(x)+\hat{\Lambda}_j'u(x)], \qquad (11)$$

where $\tilde{Q}_i(x)=(1/2)[(x^{(\hat{\Lambda}_i)}-\hat{\mu}_i)'\hat{S}_p^{-1}(x^{(\hat{\Lambda}_i)}-\hat{\mu}_i)]$, i=1,2.


A criterion for deciding whether (10) or (11) should be used is suggested in the analysis of the example in section 4 below.

### 3.1.3 Cross-validation assessment

To evaluate the performance of the rules (10) and (11), we need to compute estimates of the error rates $p_{ij}=P[\text{Assign x to } G_i|G_j]$, i.e., the probabilities of allocating an individual to $G_i$ when, in fact, it comes from $G_j$. A well established method for estimating the $p_{ij}$ is the method of cross-validation proposed by Lachenbruch and Mickey (1968). The general idea in cross-validation is determining the allocation rule using the training data deleting one observation, and then using the rule to classify the observation left out. The estimates are $\hat{p}_{ij}=a_{ij}/n_i$, where $a_{ij}$ is the number of observations from $G_i$ which have been misclassified.

For the rules (10) and (11), we need to compute $\hat{\Lambda}_{i(j)}$, the MLE of $\Lambda_i$ deleting the jth row of the data matrix $\mathcal{X}_i$, i=1,2, j=1, ...,$n_i$. This is computationally expensive and, therefore, an approximation is in order.

6

By (9), the MLE of $\Lambda_1$ is obtained by minimizing $|M_1(\Lambda_1)|$, where $M_1(\Lambda_1)=Z_1^{(\Lambda_1)'}Q_1 Z_1^{(\Lambda_1)}$. $\hat{\Lambda}_{1(j)}$ minimizes $|M_{1(j)}(\Lambda_1)|$. Let $\Lambda_{1,0}$ an initial guess for $\hat{\Lambda}_{1(j)}$ (typically $\hat{\Lambda}_1$) and let $\nabla_{1(j)}(\Lambda_1)=\partial|M_{1(j)}(\Lambda_1)|/\partial\Lambda_1$ and $H_{1(j)}(\Lambda_1)=\partial^2|M_{1(j)}(\Lambda_1)|/\partial\Lambda_1\partial\Lambda_1'$, be, respectively, the gradient vector and the hessian matrix of $|M_{1(j)}(\Lambda_1)|$. In the first-order Taylor expansion $\nabla_{1(j)}(\Lambda_1)\cong\nabla_{1(j)}(\Lambda_{1,0})+H_{1(j)}(\Lambda_{1,0})(\Lambda_1-\Lambda_{1,0})$, we can use the fact that $\nabla_{1(j)}(\hat{\Lambda}_{1(j)})=0$ to obtain a one-step approximation for $\hat{\Lambda}_{1(j)}$,

$$\hat{\Lambda}_{1(j)}^1=\Lambda_{1,0}-(H_{1(j)}(\Lambda_{1,0}))^{-1}\nabla_{1(j)}(\Lambda_{1,0}). \qquad (12)$$

Equation (12) is a multivariate extension of equation (15) in Tsai and Wu (1990).

To compute $\nabla_{1(j)}(\hat{\Lambda}_1)$ and $H_{1(j)}(\hat{\Lambda}_1)$ in (12), we define, for $i=1,2$; $j=1, \ldots,n_1$ and $k=1, \ldots,p$: (i) $a_{1jk}=x_{1jk}^{1/n_1}(\prod_{m\neq j}^{n_1}x_{1mk})^{-1/n_1(n_1-1)}$, $q_{1jk}(\lambda_{1k})=(a_{1jk})^{\lambda_{1k}-1}$; (ii) $Z_{1jk}^{*(\lambda_{1k})}=q_{1jk}(\lambda_{1k})Z_{1k}^{(\lambda_{1k})}$; and (iii) $A_{1j}=Q_1 B_{1j}Q_1$, $B_{1j}=I_{n_1}-[1-(1/n_1)]^{-1}e_{1j}e_{1j}'$, where $e_{1j}$ is the jth canonical vector of $\mathbb{R}^{n_1}$. We have the following result.

THEOREM 3.1    For $i=1,2$ and $j=1, \ldots,n_1$, define the $n_1\times p$ matrix $Z_{1j}^{*(\Lambda_1)}=(Z_{1j1}^{*(\lambda_{11})}, \ldots, Z_{1jp}^{*(\lambda_{1p})})$. We can write:

$$M_{1(j)}(\Lambda_1)=Z_{1j}^{*(\Lambda_1)}A_{1j}Z_{1j}^{*(\Lambda_1)}. \qquad (13)$$

Proof. We will use the following formula (see Atkinson (1986, p.31)): If $a,b$ are two vectors in $\mathbb{R}^n$ and $P=(p_{uv})$ is an $n\times n$ orthogonal projection matrix, we have, for $l=1, \ldots, n$,

$$a'(I_n-P)b=a_{(l)}'(I_{n-1}-P_{(l)})b_{(l)}+(1-p_{ll})^{-1}a_l^*b_l^*, \qquad (14)$$

The subscript $(l)$ means that the corresponding quantity has been computed deleting the lth row or coordinate. $a_l^*$ and $b_l^*$ are the lth coordinates of the residual vectors $a^*=(I_n-P)a$ and $b^*=(I_n-P)b$. For

$i=1,2$, $j=1$, ..., $n_1$ and $r,s=1$, ..., $p$, the $(r,s)$ element in the $p \times p$ matrix $M_{1(j)}(\Lambda_1) = (m_{1rs(j)}(\lambda_{1r}, \lambda_{1s}))$ is $m_{1rs(j)} = Z_{1r(j)}^{(\lambda_{1r})'} Q_{1(j)} Z_{1s(j)}^{(\lambda_{1s})}$. By using an analog of expression (13) in Tsai and Wu (1990) and formula (14) above, we get $m_{1rs(j)} = Z_{1jr}^{*(\lambda_{1r})} A_{1j} Z_{1js}^{*(\lambda_{1s})}$. ∎

It is important to observe that, by (13), $\hat{\Lambda}_{1(j)}$ minimizes a determinant which is structurally similar as the determinant which minimizes $\hat{\Lambda}_1$. Expressions for the partial derivatives of $|M_{1(j)}(\Lambda_1)|$ can be seen in the appendix.

## 3.2 More than two groups

When $g>2$, the construction and assessment of the discriminant rule (7) can be done following straightforward extensions of the ideas in sections 3.1.1, 3.1.2 and 3.1.3.

# 4. A PRACTICAL EXAMPLE

In DELPHI, one of the detectors of LEP, a particle collider at the European Center for Nuclear Research (CERN, Geneva, Switzerland), several variables are observed by measuring the events produced in the particle collisions. The problem is classifying, or "tagging", on the basis of the observed measurements, the event containing a quark into one of three different groups:

$G_1$: Up and Down (u+d) or Strange (s);

$G_2$: Charm (c);

$G_3$: Beauty (b).

Among the set of variables which are observed, we will consider only a special subset of 9 continuous variables $(X_1, ..., X_9)$, called Microvertex variables. We are given simulated data matrices $\mathcal{X}_1$ from the

groups $G_i$, such that $\mathcal{X}_i$ is $n_i \times 9$, with $n_1=182$, $n_2=52$, and $n_3=66$. We illustrate, with reference to this example, the discriminant techniques presented in section 3. The problem of classification has been also treated by de la Vaissière and Palma-Lopes (1989), two physicists at CERN.

The use of the rule in section 3 arises because of the strong nonnormality of the data. For example, figure 1 below shows the skewed histograms for the variable $X_4$ for the groups $G_1$, $G_2$ and $G_3$. We decided then, applying the multivariate transformation given in (5). Since, for other variables, negative data appear, we consider, for $i=1,2,3$, $j=1,\ldots,n_i$ and $k=1, \ldots,9$, the translation $x_{ijk} \longrightarrow x_{ijk} + a_k$, where $a_k = -x_{\min,k} + .5$ and $x_{\min,k}$ is the global minimum of the kth variable over the three data matrices $\mathcal{X}_1$, $\mathcal{X}_2$, $\mathcal{X}_3$, $k=1, \ldots,p$. We could have considered a shifted version of the transformation (4). Unfortunately, shifting the transformation (4) causes the maximum likelihood estimation of the transformation parameters to become a non-regular problem and not a well established procedure exists yet. However, see a recent work of Atkinson et al. (1991).

Figure 1

By maximizing the respective functions (9) for each class, we found the set of estimated transformation parameters displayed in table 1 below.

Table 1

In order to use (10) or (11) we need now to decide whether or not the dispersion matrices of the transformed data are the same. A reasonable

procedure is to test for homoscedasticity treating the transformed data matrices $x_1^{(\hat{\Lambda}_1)}$, $x_2^{(\hat{\Lambda}_2)}$ and $x_3^{(\hat{\Lambda}_3)}$ as data matrices from a multivariate normal distribution. The standard likelihood ratio test, as described for example in Mardia et al. (1979), p. 140, revealed strong evidence again homoscedasticity and, therefore, we decided to use the rule of section 3 in the form (10). It is important to remark that these procedure for choosing between (10) and (11) is not entirely rigorous because we are ignoring the effect of the sampling variability of the $\hat{\Lambda}_1$ in the null distribution of the likelihood ratio test statistic for homoscedasticity, i.e. we are treating the $\hat{\Lambda}_1$ as known constants. This procedure should be therefore taken as merely an indicative guideline for choosing between (10) and (11).

Finally, table 2 shows the estimated error rates by cross-validation. The rule classifies correctly slightly less than 70% of the events containing a beauty quark. This is remarkable because, from a physical point of view, these are the events where correct classification is more important. The relatively poor performance of the rule regarding events containing a quark of type c is not unexpected because $G_2$ is a transition class between $G_1$ and $G_3$ and, therefore, not very well differentiated from $G_1$ and $G_3$.

Table 2

## 5. CONCLUSIONS

This paper presents a new discriminant rule for the case when the groups do not have a multivariate normal distribution but can be

modelled approximately as normal after a suitable transformation. From a practical point of view, the rule presented is particularly useful in the presence of long tailed asymmetric marginal distributions which are susceptible of treatment under the framework (4.1)–(4.2). The implementation of the rule should be accompanied with a diagnostic check of normality and heteroscedasticity of the transformed data.

## APPENDIX: ELEMENTS OF $\nabla_{1(j)}(\hat{\Lambda}_1)$ AND $H_{1(j)}(\hat{\Lambda}_1)$.

For a general $n \times p$ matrix $H = (h_1, \ldots, h_p)$, we will write $H_j(u)$ for the $n \times p$ matrix obtained from $H$ by replacing its $j$th column by the $n \times 1$ vector $u$. We will also write $H_{jl}(u,v)$ for the $n \times p$ matrix obtained from $H$ by replacing its $j$th and $l$th column by, respectively, $u$ and $v$.

**THEOREM A.1** Write, for $i=1,2$, $\hat{\Lambda}_1 = (\hat{\lambda}_{11}, \ldots, \hat{\lambda}_{1p})'$ and define, for $i=1,2$, $j=1, \ldots, p$, and $k=1, \ldots, p$, the $n_1 \times 1$ vectors $Z^{*(\lambda_{1k})}_{1jk} = q_{1jk}(\lambda_{1k})Z^{(\lambda_{1k})}_{1k}$ (see section 3.1.3), $W^{*(\lambda_{1k})}_{1jk} = \partial Z^{*(\lambda_{1k})}_{1jk}/\partial \lambda_{1k}$ and $U^{*(\lambda_{1k})}_{1jk} = \partial^2 Z^{*(\lambda_{1k})}_{1jk}/\partial \lambda_{1k}^2$. We will write $Z^*_{1jk}$, $W^*_{1jk}$ and $U^*_{1jk}$ for the corresponding vectors evaluated at $\hat{\lambda}_{1k}$. If $Z^*_{1j} = Z^{*(\hat{\Lambda}_1)}_{1j}$ (see section 3.1.3), we have, for $r,s=1, \ldots, p$,

(a) $\partial |M_{1(j)}(\Lambda_1)|/\partial \lambda_{1r} \Big|_{\Lambda_1 = \hat{\Lambda}_1} = 2 |Z^{*'}_{1j,r}(W^*_{1jr})A_{1j}Z^*_{1j}|$;

(b) $\partial^2 |M_{1(j)}(\Lambda_1)|/\partial^2 \lambda_{1r} \Big|_{\Lambda_1 = \hat{\Lambda}_1} = 2[|Z^{*'}_{1j,r}(W^*_{1jr})A_{1j}Z^*_{1j,r}(W^*_{1jr})|$

$$+ |Z^{*'}_{1j}A_{1j}Z^*_{1j,r}(U^*_{1jr})|];$$

(c) $\partial^2 |M_{1(j)}(\Lambda_1)|/\partial \lambda_{1r}\partial \lambda_{1s} \Big|_{\Lambda_1 = \hat{\Lambda}_1} = 2[|Z^{*'}_{1j,rs}(W^*_{1jr},W^*_{1js})A_{1j}Z^*_{1j}|$

$$+ |Z^{*'}_{1j,r}(W^*_{1jr})A_{1j}Z^*_{1j,s}(W^*_{1js})|]. \quad (r \neq s).$$

To obtain these expressions we need the following lemma.

**LEMMA A.1** Let C, D, and B be three $n \times p$ matrices and let E be a symmetric matrix of $n \times n$. If $e_i$ represents the ith canonical vector of $\mathbb{R}^n$, we have, for $i=1, \ldots, n$, $j=1, \ldots, p$, the following differentiation formulae: a)$\partial|C'ED|/\partial c_{ij}=|C'_j(e_i)ED|$; b)$\partial|C'ED|/\partial d_{ij}=|C'ED_j(e_i)|$; and c)$\partial|B'EB|/\partial b_{ij}=2|B'_j(e_i)EB|$.

**Proof.** See lemma 3.1 in Velilla (1992). ∎

**Proof of theorem A.1.** By the chain rule and part a) of lemma A.1,

$$\partial|M_{1(j)}(\Lambda_1)|/\partial\lambda_{ir}\Big|_{\Lambda_1=\hat{\Lambda}_1}=2\sum_{i=1}^{n_1}\sum_{m=1}^{p}|Z^{*'}_{1j,m}(e_{ii})A_{1j}Z^{*}_{1j}|(\partial z^{*(\lambda_1)}_{1j1m}/\partial\lambda_{ir}\Big|_{\Lambda_1=\hat{\Lambda}_1})=$$

$$2\sum_{i=1}^{n_1}|Z^{*}_{1j,r}(w^{*}_{1j1r}e_{ii})A_{1j}Z^{*}_{1j}|=2|Z^{*'}_{1j,r}(W^{*}_{1jr})A_{1j}Z^{*}_{1j}|.$$ Expressions b) and c) are obtained in a similar way by using, respectively, parts a) and b) of lemma A.1. ∎

Finally, if from $Z^{(\lambda_{1k})}_{1k}$ (see section 3.1.1), we define $W^{(\lambda_{1k})}_{1k}=\partial Z^{(\lambda_{1k})}_{1k}/\partial\lambda_{1k}$, $U^{(\lambda_{1k})}_{1k}=\partial^2 Z^{(\lambda_{1k})}_{1k}/\partial\lambda^2_{1k}$ and $Z_{1k}$, $W_{1k}$, $U_{1k}$ stand for the corresponding functions evaluated at $\hat{\lambda}_{1k}$, it is easily seen that $W^{*}_{1jk}=q_{1jk}(\hat{\lambda}_{1k})[\log(a_{1jk})Z_{1k}+W_{1k}]$ and $U^{*}_{1jk}=q_{1jk}(\hat{\lambda}_{1k})[(\log(a_{1jk}))^2 Z_{1k} + 2\log(a_{1jk})W_{1k}+U_{1k}]$. As a conclusion, for the cross-validation assessment of the rules (10) and (11), we need the two MLE's $\hat{\Lambda}_1$ and $\hat{\Lambda}_2$, the array of constants $(a_{1jk})$ and the $n_1 \times 1$ vectors $Z_{1k}$, $W_{1k}$ and $U_{1k}$. An explicit expression for the functions $Z^{(\lambda_{1k})}_{1k}$, $W^{(\lambda_{1k})}_{1k}$ and $U^{(\lambda_{1k})}_{1k}$, can be found in Atkinson and Lawrance (1989).

## ACKNOWLEDGEMENTS

## REFERENCES

ANDREWS, D.F, GNANADESIKAN, R. & WARNER, J.L. (1971). Transformations of multivariate data. *Biometrics*, **27**, 825-840.

ATKINSON, A.C. & LAWRANCE, A.J. (1989). A comparison of asymptotically equivalent test statistics for regression transformation. *Biometrika*, **76**, 223-229.

ATKINSON, A.C., PERICCHI, L.R. & SMITH, A.F.M. (1991). Grouped likelihood for the shifted power transformation. *Journal of the Royal Statistical Society* , **B**, **53**, 473-482.

BOX. G.E.P. & COX, D.R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society* , **B**, **26**, 211-252.

DE LA VAISSIÈRE, CH. & PALMA-LOPES, S. (1989). Multidimensional analysis: a tool for b-tagging?. Unpublished manuscript.

FISHER, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, **7**, 179-188.

GNANADESIKAN, R. (1977). *Methods for Statistical Data Analysis of Multivariate Observations*. (J. Wiley: New York).

KRZANOWSKI, W. J. (1990). *Principles of Multivariate Analysis: A User's Perspective*, Paperback Edition, (Clarendon Press: Oxford, U.K).

LACHENBRUCH, P.A., & MICKEY, M.R. (1968). Estimation of error rates in discriminant analysis. *Technometrics*, **10**, 1-11.

MARDIA, K. V., KENT, J.T & BIBBY, J.M. (1979). *Multivariate Analysis*. (Academic Press: London)

SEBER, G.A.F. (1984). *Multivariate Observations*. (J. wiley: New York).

TSAI, C. L. & WU, X. (1990). Diagnostics in transformation and weighted regression. *Technometrics*, **32**, 315-322.

VELILLA, S. (1992). A note on the multivariate Box-Cox transformation to normality. Universidad Carlos III de Madrid, Working Paper 92-08. (Submitted)

# CAPTIONS FOR TABLES AND FIGURES

Table 1. Estimated transformation parameters for groups $G_1$, $G_2$ and $G_3$.

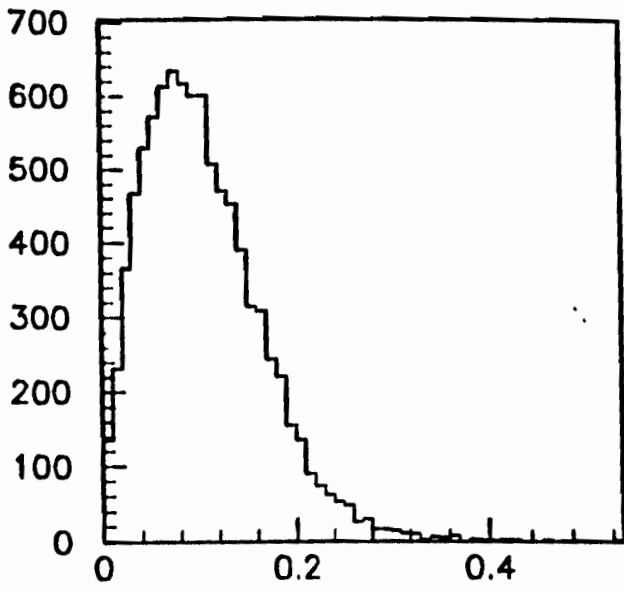Table 2. Estimated error rates by cross-validation.

Figure 1. Histograms for $X_4$ for groups $G_1$, $G_2$ and $G_3$. a) $G_1$; b) $G_2$; c) $G_3$.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| $\hat{\Lambda}'_1$ | 0.985 | 1.437 | 0.108 | 0.485 | -2.534 | -1.347 | 1.185 | 1.025 | 1.245 |
| $\hat{\Lambda}'_2$ | 0.295 | -0.313 | 0.550 | 0.558 | -1.736 | -1.043 | 1.062 | 0.853 | 0.826 |
| $\hat{\Lambda}'_3$ | -0.622 | 0.405 | 0.930 | 0.582 | -0.922 | -0.377 | 0.303 | 0.695 | 0.976 |

Table 1

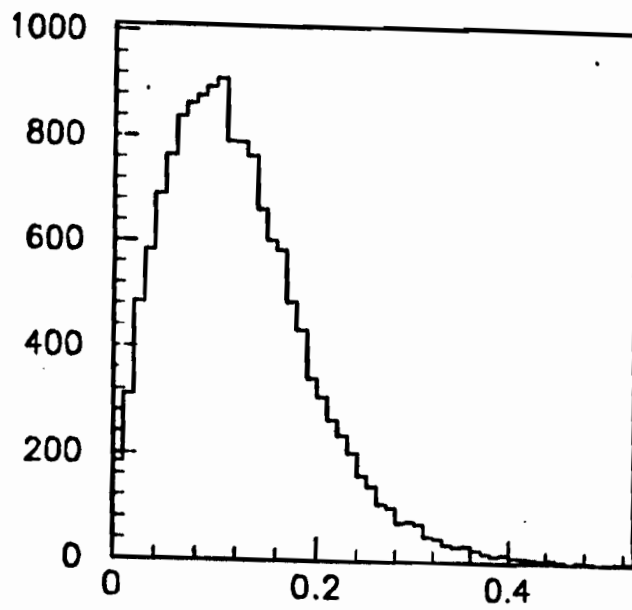| | $G_1$ | $G_2$ | $G_3$ |
|---|---|---|---|
| $G_1$ | 0.621 | 0.236 | 0.143 |
| $G_2$ | 0.481 | 0.231 | 0.288 |
| $G_3$ | 0.182 | 0.136 | 0.682 |

Table 2

15

a)

b)

c).

Figure 1