



# A two-group canonical variate analysis biplot for an optimal display of both means and cases

Niel le Roux<sup>1</sup> · Sugnet Gardner-Lubbe<sup>1</sup>

Received: 30 July 2022 / Revised: 22 March 2024 / Accepted: 27 March 2024  
© The Author(s) 2024

## Abstract

Canonical variate analysis (CVA) entails a two-sided eigenvalue decomposition. When the number of groups,  $J$ , is less than the number of variables,  $p$ , at most  $J - 1$  eigenvalues are not exactly zero. A CVA biplot is the simultaneous display of the two entities: group means as points and variables as calibrated biplot axes. It follows that with two groups the group means can be exactly represented in a one-dimensional biplot but the individual samples are approximated. We define a criterion to measure the quality of representing the individual samples in a CVA biplot. Then, for the two-group case we propose an additional dimension for constructing an optimal two-dimensional CVA biplot. The proposed novel CVA biplot maintains the exact display of group means and biplot axes, but the individual sample points satisfy the optimality criterion in a unique simultaneous display of group means, calibrated biplot axes for the variables, and within group samples. Although our primary aim is to address two-group CVA, our proposal extends immediately to an optimal three-dimensional biplot when encountering the equally important case of comparing three groups in practice.

**Keywords** Biplot · Canonical variate analysis · Classification · Data visualization · Discriminant analysis

**Mathematics Subject Classification** 62H30 · 15A21 · 93B60 · 91C20 · 97K80

## 1 Introduction

It is difficult to overrate the value of graphical displays to accompany formal statistical classification procedures (see e.g., Tukey 1975). Indeed, it is an open question if a complete assessment of group structure, including the overlap and separation of

---

✉ Niel le Roux  
njlr@sun.ac.za

<sup>1</sup> Stellenbosch University, Stellenbosch, South Africa

groups and the within groups sample behaviour is possible without the aid of suitable graphics, relying exclusively on computed statistical measures and tables. In particular, graphical procedures that are capable of displaying simultaneously the means and the individuals of different groups together with information on the relative contributions of a chosen set of classification variables are highly in demand.

Gabriel (1971) proposed the biplot for displaying simultaneously the rows and columns of a data matrix  $X : n \times p$  in a single graph. A year later Gabriel showed how to construct a canonical variate analysis (CVA) biplot (Gabriel 1972) that provides a two-dimensional graphical approximation of various groups optimally separated according to a multidimensional CVA criterion. This CVA biplot became popular among statisticians and practitioners applying linear discriminant analysis (LDA) and CVA in various fields. Gittins (1985) can be consulted for an overview of CVA. Distances in this classical Gabriel CVA biplot are interpreted in terms of inner products between vectors. Such interpretations are not as straightforward as distances in an ordinary scatter plot. The unified biplot methodology proposed by Gower (1995) and extensively discussed in the monograph by Gower and Hand (1996) allows the CVA biplot to be regarded as a multivariate extension of an ordinary scatter plot. The concept of inter-sample distance is central in this approach, while information on the classifier variables is added to the graph of means and individual sample points in the form of biplot axes—an axis calibrated in its original units for each (classifier) variable. The perspective of Gower and Hand (1996) inspires the biplot methodology discussed in this paper. Gower et al. (2011) discuss the one-dimensional CVA biplot for use in two-group studies in some detail. Although the latter authors provide several enhancements to this one-dimensional CVA biplot they failed to address the challenge of constructing an optimal two-dimensional biplot for the two-group case.

Classification procedures in practice are often confronted with the problem of optimally distinguishing between two groups. The aim of several procedures used in multidimensional classification is to separate the group means optimally. This is the aim of CVA and the closely related procedure of LDA. These techniques involve the transformation of the original observations into a so-called canonical space. Flury (1997), among others, proves that in the case of  $J$  groups, only the first  $K = \min(p, J - 1)$  elements of the group mean vectors differ in the canonical space. This induces some pitfalls when routinely constructing CVA biplots in the case where  $p > J - 1$  with  $J$  equaling two or three. Since in the case of two groups the group means can be exactly represented on a line, the associated CVA biplot becomes a one-dimensional plot with all  $p$  biplot axes representing the different variables on top of each other on the line extending through the representations of the two group means. All  $n$  sample points also fall on this line. The theoretical basis for constructing this line is provided by the eigenanalysis of a two-sided eigenvalue problem (see e.g., Gower and Hand (1996) and Gower et al. (2011)). In the two-sample case this eigenanalysis involves only one non-zero (positive) eigenvalue together with  $p - 1$  zero eigenvalues. Therefore, only the eigenvector associated with the single non-zero eigenvalue is uniquely defined. This eigenvector provides the scaffolding for constructing a one-dimensional CVA biplot. Although a two-dimensional biplot can be constructed using two eigenvectors, the second eigenvector will not be uniquely defined, as is the case when  $J > 2$ , unless some precautions are taken. A similar situation arises in the case

of three groups: the first two eigenvectors are then uniquely defined but not the third. Why do we want extra scaffolding dimensions for constructing CVA biplots when  $p > J - 1$ ? There is a real advantage when we have a two-group or three-group classification problem: not only are the group means then represented exactly but also an improvement in the approximations of the individual sample points is accomplished. Therefore, in this paper, we first define a measure of how well any mean or individual sample point is represented in a CVA biplot. Then we show, in the case of two groups, how to construct a uniquely defined two-dimensional CVA biplot such that both the means and individual samples are optimally represented together with the variables in the form of calibrated biplot axes. In addition, we will show how this process can be extended directly for constructing a uniquely defined optimal three-dimensional biplot when three groups are considered.

The paper is organized as follows: in the next section, we begin with some historical background of discriminant analysis. After that, we briefly review the basic concepts and theory of CVA biplot methodology according to the perspective of Gower and Hand (1996). This is followed by a section describing the geometry of the graphical representation of the class means and the sample points about them. In section 4 it is discussed why the two-group CVA biplot deserves special attention. We then put forward proposals for one-dimensional CVA biplots as well as a two-dimensional CVA biplot such that the group means in a two-group classification problem are exactly represented together with an optimal representation of the individual sample points. We also show how to generalize this procedure to construct a three-dimensional biplot with similar properties for use when  $J = 3 < p$ . In section 6 we briefly discussed an alternative unique 2D biplot, which is based on the Bhattacharyya distance. The theoretical results are illustrated in section 7 where we provide examples, covering one- and two-dimensional CVA biplots. Some conclusions are considered in section 8.

## 2 A brief review of canonical variate analysis

### 2.1 Two-group discriminant analysis: Fisher's LDA

Two-group discriminant analysis considers two populations (groups)  $G_1$  and  $G_2$ . An observation  $\mathbf{x}$  of  $\mathbf{X}^T = (X_1, X_2, \dots, X_p)$  is to be allocated to one of these populations. It is assumed that the density  $f_i(\mathbf{x})$  of  $\mathbf{X}$  for  $i = (1; 2)$  is known with expected value  $\boldsymbol{\mu}_i : p \times 1$  and covariance matrix  $\boldsymbol{\Sigma}_i : p \times p$ , respectively. Let the prior probability of an unknown  $\mathbf{x}$  belonging to  $G_i$  be given by

$$p_1 = P(G_1) \text{ and } p_2 = P(G_2), \text{ respectively, with } p_i > 0 \text{ and } p_1 + p_2 = 1.$$

Define

$$\begin{aligned} \boldsymbol{\mu} &= E(\mathbf{X}) = p_1\boldsymbol{\mu}_1 + p_2\boldsymbol{\mu}_2, \\ \mathbf{T} &= E\left[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\right], \\ \mathbf{B} &= p_1(\boldsymbol{\mu}_1 - \boldsymbol{\mu})(\boldsymbol{\mu}_1 - \boldsymbol{\mu})^T + p_2(\boldsymbol{\mu}_2 - \boldsymbol{\mu})(\boldsymbol{\mu}_2 - \boldsymbol{\mu})^T, \text{ and} \\ \mathbf{W} &= p_1\boldsymbol{\Sigma}_1 + p_2\boldsymbol{\Sigma}_2. \end{aligned}$$

Under the assumption that  $\Sigma_1 = \Sigma_2 = \Sigma$ , (say) it follows that  $W = \Sigma$ . The Fisher LDA (Fisher 1936) searches for the linear function

$$Y = m^T X$$

with  $E(Y) = m^T E(X) = m^T \mu$ ,  $E(Y|G_i) = m^T \mu_i$  and  $var(Y) = m^T W m$  to maximize

$$\frac{p_1(m^T \mu_1 - m^T \mu)^2 + p_2(m^T \mu_2 - m^T \mu)^2}{m^T W m} = \frac{m^T B m}{m^T W m}. \tag{1}$$

The maximum is obtained from the eigenequation

$$(W^{-1} B)m = m \Lambda.$$

Pre-multiplying the above equation with  $W^{1/2}$  leads to

$$(W^{-1/2} B W^{-1/2})(W^{1/2} m) = (W^{1/2} m) \Lambda,$$

so that  $l = W^{1/2} m$  is the eigenvector of  $W^{-1/2} B W^{-1/2}$  associated with the largest eigenvalue  $\lambda_1$  and  $m = W^{-1/2} l$ .

Since  $rank(B) = 1 = rank(W^{-1} B) = rank(W^{-1/2} B W^{-1/2})$ , it follows that

$$B M = W M \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix}. \tag{2}$$

If we choose  $L$  to be an orthogonal matrix in  $(W^{-1/2} B W^{-1/2})L = L \Lambda$ , then

$$L^T L = I = M^T W M$$

and

$$M^T B M = M^T W M \Lambda = \Lambda.$$

The transformation

$$Y_k = m_k^T X$$

for  $k = 1, 2, \dots, p$  is termed a transformation into the canonical space with  $Y_1$  the first canonical variate, where

$$E(Y_1|G_i) = m_1^T \mu_i \text{ and } Var(Y_1|G_i) = Var(Y_1) = 1 \text{ for } i = 1; 2.$$

After properly scaled, the solution  $\mathbf{m}_1$  maximizing (1) can be written as

$$\mathbf{W}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \tag{3}$$

which is known as Fisher’s linear discriminant function (LDF). The maximum of (1) is given by the squared Mahalanobis distance, namely

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \mathbf{W}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

For  $k = 2, 3, \dots, p$  we have  $Var(Y_k|G_i) = 1$  and

$$\mathbf{m}_k^T \mathbf{B} \mathbf{m}_k = 0. \tag{4}$$

Since

$$\begin{aligned} \mathbf{B} &= p_1(\boldsymbol{\mu}_1 - \boldsymbol{\mu})(\boldsymbol{\mu}_1 - \boldsymbol{\mu})^T + p_2(\boldsymbol{\mu}_2 - \boldsymbol{\mu})(\boldsymbol{\mu}_2 - \boldsymbol{\mu})^T \\ &= (p_1 p_2^2 + p_2 p_1^2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T, \end{aligned}$$

it follows from (4) that  $\mathbf{m}_k^T(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = 0$  and all differences vanish between the groups for the second and further canonical variates.

Rao (1948) extends the Fisher’s LDA procedure to  $J > 2$  groups by deriving  $J - 1$  linear discriminant functions of the form (3), but in this paper we are primarily interested in the case  $J = 2$ .

### 2.2 Bayes linear and quadratic classifiers

Under the assumption of multivariate normal distributions with different covariance matrices  $\boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$ , the Bayes quadratic classifier (see e.g., Hastie et al. 2001) is given by

$$\begin{aligned} &\frac{1}{2}(X - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1}(X - \boldsymbol{\mu}_1) - \frac{1}{2}(X - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1}(X - \boldsymbol{\mu}_2) \\ &+ \frac{1}{2} \log\left(\frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_2|}\right) \stackrel{\geq}{\leq} \log\left(\frac{p_1}{p_2}\right). \end{aligned}$$

If  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$ , (say) we have the Bayes linear classifier

$$\frac{1}{2}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} X + \frac{1}{2}(\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2) \stackrel{\geq}{\leq} \log\left(\frac{p_1}{p_2}\right).$$

It is clear that if  $p_1 = p_2$ , the Bayes linear classifier is equivalent to Fisher’s LDF.

In general, the Bhattacharyya distance measures the similarity of two probability distributions. When the distributions concerned are  $N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  and  $N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$  it is

given by (see e.g., Fukunaga 1990; Hennig 2004):

$$D_{Bhat} = \frac{1}{8}(\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \left( \frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right)^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \frac{1}{2} \log \left( \frac{\left| \frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right|}{\sqrt{|\boldsymbol{\Sigma}_1| |\boldsymbol{\Sigma}_2|}} \right). \tag{5}$$

It can be shown that the sample form of (5) provides an upper bound for the Bayes error (see e.g., McLachlan 1992). Furthermore, when  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$  then (5) is proportional to a squared Mahalanobis distance. For this case, Fukunaga (1990) shows that  $D_{Bhat}$  is maximized by  $Y = \mathbf{m}^T X$  where  $\mathbf{Bm} = \lambda_1 \mathbf{Wm}$ , so that maximization is achieved when  $\mathbf{m}$  is taken as (3).

### 2.3 Two-group discrimination where group means and or group covariance matrices may differ

While LDA discussed above allows only for groups to differ with respect to the means, Fukunaga (1990) looks for a linear transformation  $Y = X\mathbf{M}$  to separate groups with respect to their means or covariance matrices. With the notation of the subsection above, write

$$\begin{aligned} \mathbf{W} &= p_1 \boldsymbol{\Sigma}_1 + p_2 \boldsymbol{\Sigma}_2 = p_1 E \left[ (X - \boldsymbol{\mu}_1)(X - \boldsymbol{\mu}_1)^T | G_1 \right] \\ &+ p_2 E \left[ (X - \boldsymbol{\mu}_2)(X - \boldsymbol{\mu}_2)^T | G_2 \right]. \end{aligned}$$

Fukunaga (1990) provides four criteria for class separability:

$J_1 = tr(\mathbf{S}_2^{-1} \mathbf{S}_1)$ ;  $J_2 = \log |\mathbf{S}_2^{-1} \mathbf{S}_1|$ ;  $J_3 = tr(\mathbf{S}_1) - \lambda(tr(\mathbf{S}_2) - c)$ , where  $\lambda$  is a Lagrange multiplier,  $c$  is a constant, and  $J_4 = tr(\mathbf{S}_1)/tr(\mathbf{S}_2)$ , where  $\mathbf{S}_1$  and  $\mathbf{S}_2$ , respectively, are one of  $\mathbf{T}$ ,  $\mathbf{B}$ , or  $\mathbf{W}$ .

Let  $J_i(k)$  indicate the criterion with  $\mathbf{S}_1 = \mathbf{B}$  and  $\mathbf{S}_2 = \mathbf{W}$ , and where  $Y = X\mathbf{M}$  with  $\mathbf{M} : p \times k$ . Since we are restricted to linear transformations, optimization of  $J_1(1)$  is equivalent to Fisher discriminant analysis with (1) and  $J_1(1) = \lambda_1$  with no other dimension contributing to the value of  $J_1$ . Fukunaga (1990) further shows that criterion  $J_1$  gives the same optimum transformation for other combinations of  $\mathbf{B}$ ,  $\mathbf{W}$ , and  $\mathbf{T}$  for  $\mathbf{S}_1$  and  $\mathbf{S}_2$  and also for optimizing  $J_2$ .

When  $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \boldsymbol{\mu}$ , (5) becomes

$$\begin{aligned} D_{Bhat} &= \frac{1}{2} \log \left( \frac{\left| \frac{\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2}{2} \right|}{\sqrt{|\boldsymbol{\Sigma}_1| |\boldsymbol{\Sigma}_2|}} \right) \\ &= \frac{1}{4} \left[ \log \left( |\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_2 + 2\mathbf{I}_p| \right) - p \log(4) \right]. \end{aligned} \tag{6}$$

If more than a single dimension is needed we have that

$$J_1 = \log |(\mathbf{M}^T \boldsymbol{\Sigma}_2 \mathbf{M})^{-1} (\mathbf{M}^T \boldsymbol{\Sigma}_1 \mathbf{M})| + (\mathbf{M}^T \boldsymbol{\Sigma}_1 \mathbf{M})^{-1} (\mathbf{M}^T \boldsymbol{\Sigma}_2 \mathbf{M}) + 2\mathbf{I}_k|$$

is maximized by  $Y = XM$  with  $M : p \times k$ , where

$$(\Sigma_2^{-1} \Sigma_1)M = M \left( (M^T \Sigma_2 M)^{-1} (M^T \Sigma_1 M) \right)$$

and

$$(\Sigma_1^{-1} \Sigma_2)M = M \left( (M^T \Sigma_1 M)^{-1} (M^T \Sigma_2 M) \right).$$

Thus,  $M$  must contain the eigenvectors of both  $\Sigma_2^{-1} \Sigma_1$  and  $\Sigma_1^{-1} \Sigma_2$ . However, they have the same eigenvectors since  $\Sigma_2^{-1} \Sigma_1 = (\Sigma_1^{-1} \Sigma_2)^{-1}$  and we have that

$$(M^T \Sigma_2 M)^{-1} (M^T \Sigma_1 M) = \Lambda$$

and

$$(M^T \Sigma_1 M)^{-1} (M^T \Sigma_2 M) = \Lambda^{-1},$$

so that  $(\Sigma_2^{-1} \Sigma_1)M = M\Lambda$  and  $(\Sigma_1^{-1} \Sigma_2)M = M\Lambda^{-1}$ .

Therefore, if  $k > 1$  dimensions are needed the  $k$  eigenvectors are chosen, which correspond to the  $k$  largest values for  $J$ , i.e., corresponding to the largest values  $\lambda_i + \frac{1}{\lambda_i} + 2$ .

### 3 Discriminant analysis with sampled data

In practice, discriminant analysis is usually performed using sampled data. This necessitates the substitution of population parameters with sample estimates in the formulas introduced in section 2. The plug-in principle and maximum likelihood method are popular methods for this.

It is well known that LDA often outperforms quadratic discriminant analysis (QDA) even when the assumption of equal covariance matrices is violated (see e.g., Flury et al. 1997; McLachlan 1992). This can be attributed to the large number of parameters that have to be estimated in QDA with over-parameterization inducing a loss of power. This contributes to the popularity of LDA among practitioners and so in the rest of the paper, our focus will be on LDA.

Consider the data matrix  $X : n \times p$  centered such that  $\mathbf{1}^T X = \mathbf{0}^T$ . The data contained in  $X$  consists of  $p$  measurements made for each of the  $J$  groups. The group sizes are  $n_1, n_2, \dots, n_J$ , respectively, such that  $\sum_{i=1}^J n_i = n$ . Let  $N_g = \text{diag}(n_1, n_2, \dots, n_J)$ , so that a matrix of group means can be calculated as

$$\bar{X} : J \times p = N_g^{-1} G^T X = (G^T G)^{-1} G^T X, \tag{7}$$

where  $G : n \times J$  denotes an indicator matrix defining the  $J$  groups.

Let  $\mathcal{V}(X^T)$  denote the vector space generated by the columns of  $X^T$ . We assume this vector space of  $p$ -vectors to be of dimension  $p$ . Since each row of  $\bar{X}$  is a linear combination of the rows of  $X$  it follows that  $\bar{X}^T \in \mathcal{V}(X^T)$ .

Define:

1.  $S_B : p \times p$ , as the between-group matrix of squares and products:  $S_B = \bar{X}^T N_g \bar{X} = X^T G(G^T G)^{-1} G^T X$  and
2.  $S_W : p \times p$  as the within-group matrix of squares and products:  $S_W = X^T X - \bar{X}^T N_g \bar{X} = X^T (I - G(G^T G)^{-1} G^T) X$ .

Generally,  $rank(S_W) = p$  while  $rank(S_B) = \min(J - 1, p)$ .

The two-sided eigenvalue problem

$$(S_B)B = (S_W)BA \tag{8}$$

provides the solution  $b_1$  to the CVA criterion

$$\underset{b}{\text{maximize}} \left( \frac{b^T (S_B) b}{b^T (S_W) b} \right), \text{ subject to } b^T (S_W) b = 1. \tag{9}$$

In the above, the diagonal matrix  $A : p \times p$  contains the eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ , where  $\lambda_J = \lambda_{J+1} = \dots = \lambda_p = 0$  if  $J < p + 1$ .

The matrix  $B = [b_1, b_2, \dots, b_p]$  contains all  $p$  solutions to the two-sided eigenvalue problem. Only the solution  $b_1$  is optimal for the CVA criterion (9). The matrix  $B : p \times p$  is non-singular with  $B^{-1} = B^T (S_W)$ , while the columns of  $B$  are orthogonal in the metric  $S_W$  because of the constraints  $B^T (S_W) B = I$ .

Canonical variates are defined by the transformation  $y^T = x^T B$ , where  $x$  is any  $p$ -vector belonging to  $\mathcal{V}(X^T)$ . The centered data matrix itself is transformed to the canonical variate values matrix  $Y : n \times p$  through the one-to-one (canonical) transformation

$$Y : n \times p = XB. \tag{10}$$

The transformation (10) implies a transformation of  $\mathcal{V}(X^T)$  to  $\mathcal{V}(Y^T)$ , the canonical space, of dimension  $p$  since  $rank(Y) = rank(X)$ . Furthermore, (10) implies that

$$\bar{X}B = \bar{Y} : J \times p. \tag{11}$$

We will call  $\bar{Y}$  the canonical means matrix. It follows that the columns of  $\bar{Y}^T$  generate a subspace of dimension  $\min(J - 1, p)$  of  $\mathcal{V}(Y^T)$ . This subspace is denoted by  $\mathcal{V}(\bar{Y}^T)$ .



#### 4 Biplot display of the canonical means matrix $\bar{Y}$ and the canonical variates values matrix $Y$

An  $r$ -dimensional canonical variate analysis (CVA) plot is constructed by taking the first  $r$  canonical variates, associated with

$$B_r = [b_1, b_2, \dots, b_r], \text{ where } B = [B_r, b_{r+1}, \dots, b_p] \tag{12}$$

to provide the coordinates (or scaffolding) for representing the  $J$  canonical group means contained in (11) as points in  $r$  dimensions. If this plot is equipped with  $p$  linear axes to represent the original  $p$  variables, a CVA biplot is obtained. Each of these biplot axes is determined by a vector, which induces also a graduation on it. Gower and Hand (1996) consider two types of CVA biplots, each one characterized by its system of  $p$  linear axes, its aim and its corresponding geometry:

- The *interpolation* biplot, which has the aim of placing on the plot the image  $(y_1, y_2, \dots, y_r)$  of any new point  $x \in \mathcal{V}(X^T)$ .
- The *prediction* biplot, which has the aim of estimating the point  $x$  (i.e., the set of variable values) having as an image a given point  $(y_1, y_2, \dots, y_r)$  in the plot.

Once the CVA biplot for representing the group means is constructed, all transformed samples contained in  $Y$  can be interpolated into the biplot. Thus both the canonical means and the transformed samples  $Y$  can be displayed in a CVA biplot in an  $r$ -dimensional subspace of  $\mathcal{V}(\bar{Y}^T)$ , (with  $r \leq \min(J - 1, p)$ ). Typically, an  $r$  of two or three will be chosen to construct this subspace that we will call the biplot space.

Gower and Hand (1996) show the above processes of prediction and interpolation to be based on the following: A sample  $x : p \times 1$  can be interpolated into  $\mathcal{V}(Y^T)$  by

$$y : p \times 1 = B^T x, \text{ i.e., } y^T = x^T B = \sum_{k=1}^p (x_k e_k^T) B.$$

The representation of  $x$  in the biplot space is given by

$$z^T : 1 \times r = x^T B_r = \sum_{k=1}^p (x_k e_k^T) B_r, \tag{13}$$

where  $B_r$  is defined in (12).

Prediction is the inverse of interpolation and since  $B$  is non-singular it follows by inverting the above formula for interpolation that

$x^T = y^T B^{-1}$ . The matrix  $B^{-1}$  can be partitioned into

$$B^{-1} = \begin{bmatrix} B^{(r)} : r \times p \\ B^{(2)} : p - r \times p \end{bmatrix}. \tag{14}$$

The predicted value for the  $k$ th variable can be written as  $\hat{x}_k = \mathbf{z}^T \mathbf{B}^{(r)} \mathbf{e}_k$  and therefore, the predicted value for  $\mathbf{x}^T$  is

$$\begin{aligned} \hat{\mathbf{x}}^T &= \mathbf{z}^T \mathbf{B}^{(r)} [\mathbf{e}_1, \dots, \mathbf{e}_p] \\ &= \mathbf{z}^T \mathbf{B}^{(r)} \\ &= \mathbf{x}^T \mathbf{B}_r \mathbf{B}^{(r)}. \end{aligned} \tag{15}$$

It follows from (15) that

$$\hat{\mathbf{X}} = \mathbf{X} \mathbf{B}_r \mathbf{B}^{(r)}, \tag{16}$$

and in addition

$$\hat{\bar{\mathbf{X}}} = \bar{\mathbf{X}} \mathbf{B}_r \mathbf{B}^{(r)}. \tag{17}$$

Since in the CVA biplot described above, the samples are interpolated into the biplot constructed for the canonical means, it is expected that the class means will be better represented than the canonical variate values i.e., the rows of  $\mathbf{Y}$ . What is needed then, are measures of fit for use in CVA.

#### 4.1 Measures of fit for use in CVA

From the identity

$$\mathbf{X} = \mathbf{G}(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T \mathbf{X} + (\mathbf{I} - \mathbf{G}(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T) \mathbf{X},$$

we have that

$$\mathbf{X}^T \mathbf{X} = \mathbf{X}^T \mathbf{Q} \mathbf{X} + \mathbf{X}^T (\mathbf{I} - \mathbf{Q}) \mathbf{X},$$

where the matrix  $\mathbf{Q} : n \times n = \mathbf{G}(\mathbf{G}^T \mathbf{G})^{-1} \mathbf{G}^T$  is positive semi-definite, symmetric and idempotent.

This between-and-within group structure is of interest for:

- Assigning a given sample to its most appropriate group.
- Relating the groups to one another.

A study of the relationships among the groups encourages the use of low-dimensional approximations for visualization, including CVA biplots.

#### 4.2 Measures of fit for CVA biplots: recovering the canonical group means

The orthogonal partitioning

$$\mathbf{B}^T \bar{\mathbf{X}}^T N_g \bar{\mathbf{X}} \mathbf{B} = \mathbf{B}^T (\hat{\bar{\mathbf{X}}})^T N_g \hat{\bar{\mathbf{X}}} \mathbf{B} + \mathbf{B}^T (\bar{\mathbf{X}} - \hat{\bar{\mathbf{X}}})^T N_g (\bar{\mathbf{X}} - \hat{\bar{\mathbf{X}}}) \mathbf{B},$$

(see Gardner-Lubbe et al. 2008) allows an overall measure of how well the group means are represented in the CVA biplot, namely

$$\text{Overall quality} = \frac{\text{tr}\left(\mathbf{B}^T (\hat{\mathbf{X}})^T N_g \hat{\mathbf{X}} \mathbf{B}\right)}{\text{tr}\left(\mathbf{B}^T \bar{\mathbf{X}}^T N_g \bar{\mathbf{X}} \mathbf{B}\right)} = \frac{\sum_{k=1}^r \lambda_k}{\sum_{k=1}^p \lambda_k}.$$

In the orthogonal partitioning above, the matrix  $\mathbf{B}$  can be eliminated to define *axis predictivities* as the diagonal elements of the matrix

$$\mathbf{\Pi} : p \times p = \text{diag}\left((\hat{\mathbf{X}})^T N_g \hat{\mathbf{X}}\right) \left[\text{diag}(\bar{\mathbf{X}}^T N_g \bar{\mathbf{X}})\right]^{-1}.$$

Each axis predictivity is a measure of how well the values for the group means can be determined from the CVA biplot for the variable associated with that particular biplot axis. We note that the overall quality is a weighted mean of the individual axis predictivities.

### 4.3 Measures of fit for CVA biplots: recovering the individual samples

We also need predictivities for individual samples corrected for class means i.e., for  $(\mathbf{I} - \mathbf{Q})\mathbf{X}$ . So, our starting point becomes the decomposition

$$(\mathbf{I} - \mathbf{Q})\mathbf{X}\mathbf{B} = (\mathbf{I} - \mathbf{Q})\hat{\mathbf{X}}\mathbf{B} + (\mathbf{I} - \mathbf{Q})(\mathbf{X} - \hat{\mathbf{X}})\mathbf{B},$$

where  $\hat{\mathbf{X}}$  is defined in (16). Then the following orthogonal decompositions (see Gower et al. 2011) hold:

1. *Type A*

$$\begin{aligned} \mathbf{B}^T \mathbf{X}^T (\mathbf{I} - \mathbf{Q})\mathbf{X}\mathbf{B} &= \mathbf{B}^T \hat{\mathbf{X}}^T (\mathbf{I} - \mathbf{Q})\hat{\mathbf{X}}\mathbf{B} \\ &\quad + \mathbf{B}^T (\mathbf{X} - \hat{\mathbf{X}})^T (\mathbf{I} - \mathbf{Q})(\mathbf{X} - \hat{\mathbf{X}})\mathbf{B}. \end{aligned}$$

2. *Type B*

$$\begin{aligned} (\mathbf{I} - \mathbf{Q})\mathbf{X}\mathbf{B}\mathbf{B}^T \mathbf{X}^T (\mathbf{I} - \mathbf{Q}) &= (\mathbf{I} - \mathbf{Q})\hat{\mathbf{X}}\mathbf{B}\mathbf{B}^T \hat{\mathbf{X}}^T (\mathbf{I} - \mathbf{Q}) \\ &\quad + (\mathbf{I} - \mathbf{Q})(\mathbf{X} - \hat{\mathbf{X}})\mathbf{B}\mathbf{B}^T (\mathbf{X} - \hat{\mathbf{X}})^T (\mathbf{I} - \mathbf{Q}). \end{aligned}$$

While the class means are exactly represented in a subspace of dimension  $\min(J - 1, p)$  of the canonical space, this is generally not true for the individual samples. From Type B orthogonality within-group sample predictivities can be defined as the diagonal elements of

$$\Psi_W : n \times n = \text{diag}\left((\mathbf{I} - \mathbf{Q})\hat{\mathbf{X}}S_W^{-1}\hat{\mathbf{X}}^T(\mathbf{I} - \mathbf{Q})\right)\left[\text{diag}\left((\mathbf{I} - \mathbf{Q})\mathbf{X}S_W^{-1}\mathbf{X}^T(\mathbf{I} - \mathbf{Q})\right)\right]^{-1}$$

#### 4.4 CVA biplots for $J = 2$ groups

When  $J = 2$  it now follows that:

- The underlying two-sided eigenequation has one non-zero eigenvalue and  $p - 1$  zero eigenvalues.
- All  $p$  class means are exactly represented in a single dimension.
- This single dimension contains the  $p$  biplot axes (each with predictivity 100% for recovering the group means) on top of each other.
- Overall quality of representing group means is 100%.
- The one-dimensional CVA biplot is optimal for representing groups irrespective of the number of variables  $p$ .
- The samples are not exactly represented in the one-dimensional CVA biplot.

Our challenge is now to add another dimension for improving recovering of sample information without changing optimality properties already available for the group means. We address this challenge by considering the orthonormal complement in the canonical space of the subspace containing the two group means. Therefore, we consider eigenvectors associated with the zero eigenvalues. These eigenvectors have no natural ordering associated with them. So, any one of these eigenvectors or even any linear combination of them has the same claim to be used as a second scaffolding axis. Hence, our aim is to find the linear combination that satisfies some optimality criterion using Type A and Type B orthogonality for natural candidates.

#### 4.5 Optimal two-dimensional CVA biplot for $J = 2$ groups: optimality criterion based on Type B orthogonality

Consider minimizing

$$\text{sum}\left\{\text{diag}\left((\mathbf{I} - \mathbf{Q})(\mathbf{X} - \hat{\mathbf{X}})\mathbf{B}\mathbf{B}^T(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{I} - \mathbf{Q})\right)\right\}.$$

Now,

$$\begin{aligned} (\mathbf{I} - \mathbf{Q})(\mathbf{X} - \hat{\mathbf{X}}) &= (\mathbf{I} - \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T)(\mathbf{X} - \hat{\mathbf{X}}) \\ &= (\mathbf{X} - \hat{\mathbf{X}}) - \mathbf{G}(\mathbf{G}^T\mathbf{G})^{-1}\mathbf{G}^T(\mathbf{X} - \hat{\mathbf{X}}) \\ &= (\mathbf{X} - \hat{\mathbf{X}}) - \mathbf{G}\bar{\mathbf{X}} + \mathbf{G}\hat{\mathbf{X}}. \end{aligned}$$

We have  $\mathbf{G}\hat{\bar{\mathbf{X}}} = \mathbf{G}\bar{\mathbf{X}}$  if the eigenvector associated with  $\lambda > 0$  is chosen. Therefore, maximizing the sum of within-group sample predictivities becomes equivalent

to minimizing  $sum \left\{ diag \left( (X - \hat{X}) B B^T (X - \hat{X})^T \right) \right\}$ . However, this sum remains constant for any additional eigenvector. This is not surprising because in the canonical space the constraint  $B^T S_W B = I$ , implies constant variation in all dimensions, resulting in Type B orthogonality not useful to define a criterion for an optimal two-dimensional biplot.

**4.6 Optimal two-dimensional CVA biplot for  $J = 2$  groups: optimality criterion based on Type A orthogonality**

Since the matrix  $B$  is non-singular it can be eliminated from the equation defining Type A orthogonality. As before,  $(I - Q)(X - \hat{X}) = (X - \hat{X})$ . Therefore, the proposed optimality criterion for constructing an optimal two-dimensional CVA biplot for two groups is the total squared reconstruction error for samples:

$$TSRES = tr \left\{ (X - \hat{X})(X - \hat{X})^T \right\}. \tag{18}$$

A similar measure of the goodness of approximations of the means can be defined as the total squared reconstruction error for means:

$$TSREM = tr \left\{ (\bar{X} - \hat{\bar{X}})(\bar{X} - \hat{\bar{X}})^T \right\}. \tag{19}$$

**5 Optimal CVA biplots when the number of groups is less than or equal to the number of variables**

From now on we consider the case where  $J < p + 1$ . Then it follows that the canonical means matrix  $\bar{Y}$  is of the form

$$\begin{bmatrix} \bar{y}_1^T \\ \dots \\ \bar{y}_J^T \end{bmatrix} = \begin{bmatrix} k_{11} & \dots & k_{1(J-1)} & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ k_{J1} & \dots & k_{J(J-1)1} & 0 & \dots & 0 \end{bmatrix}, \tag{20}$$

(see e.g., Flury, 1997, p. 491).

Hence, the canonical transformation optimally separates the  $p$ -vectors  $\bar{y}_1^T, \bar{y}_2^T, \dots, \bar{y}_J^T$  in  $J - 1$  dimensions while all differences among them vanish in dimensions  $J, J + 1, \dots, p$ . It follows that  $\mathcal{V}(\bar{Y}^T)$  is now of dimension  $J - 1$  and it is thus possible to consider a biplot space of dimension  $r > J - 1$ . The resulting  $r$ -dimensional CVA biplot is not uniquely defined because the first  $J - 1$  columns of  $B : p \times p$  appended with any set of  $r - J + 1$  of its remaining columns will result in a biplot where the canonical means are exactly represented. Therefore,  $\bar{x}_j^T B = [k_{j1} \dots k_{j(J-1)} 0 \dots 0]$  for  $j = 1, 2, \dots, J$  with the predicted value for  $\bar{x}_j$

given by the  $j$ th row,  $\hat{\mathbf{x}}_j^T$ , of (17). It follows that

$$\begin{aligned} \hat{\mathbf{x}}_j^T &= \bar{\mathbf{x}}_j^T \mathbf{B}_r \mathbf{B}^{(r)} \\ &= \bar{\mathbf{x}}_j^T [b_1 \ b_2 \ \dots \ b_{J-1} \ b_J \ \dots \ b_r] \begin{bmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \dots \\ \mathbf{b}^{(r)} \end{bmatrix}, \text{ where } \begin{bmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \dots \\ \mathbf{b}^{(r)} \end{bmatrix} = \mathbf{B}^{(r)} \\ &= \bar{\mathbf{x}}_j^T [b_1 \ b_2 \ \dots \ b_{J-1} \ b_J \ \dots \ \mathbf{0}] \begin{bmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \dots \\ \mathbf{b}^{(p)} \end{bmatrix}, \text{ with } \begin{bmatrix} \mathbf{b}^{(1)} \\ \mathbf{b}^{(2)} \\ \dots \\ \mathbf{b}^{(p)} \end{bmatrix} = \mathbf{B}^{-1} \\ &= \bar{\mathbf{x}}_j^T \text{ for } j = 1, 2, \dots, J. \end{aligned}$$

Therefore, in the above  $r$ -dimensional CVA biplot the canonical means are exactly represented, resulting in  $TSREM = 0$ . For a sample  $\mathbf{x}_i^T$ , we have the  $i$ th row of  $\mathbf{X}$ ,

$$\mathbf{x}_i^T = \mathbf{x}_i^T \mathbf{B}_r \mathbf{B}^{(r)} \neq \mathbf{x}_i^T \mathbf{B} \mathbf{B}^{-1},$$

since in general  $\mathbf{x}_i^T \mathbf{b}_j$  is not zero for all  $j = J, J + 1, \dots, p$  so that  $TSRES > 0$ .

This leaves us with the challenge to construct scaffolding axes  $j = J, J + 1, \dots, r$  in addition to those contained in  $\mathbf{B}_{J-1}$  so that  $TSRES$  is minimized without sacrificing what we already have for the class means in  $J - 1$  dimensions.

Possible candidates for the additional scaffolding axes are any  $r - J + 1$  of the vectors  $\mathbf{b}_J, \mathbf{b}_{J+1}, \dots, \mathbf{b}_p$ . All these vectors are associated with the zero eigenvalues (diagonal elements of  $\mathbf{\Lambda}$ ). Therefore, there is no natural ordering of them as is the case with the  $J - 1$  eigenvectors associated with the non-zero eigenvalues. Furthermore, since  $\bar{\mathbf{X}}\mathbf{b}_i = \mathbf{0}$  for  $i = J, J + 1, \dots, p$  it follows that  $\bar{\mathbf{X}}\mathbf{d} = \mathbf{0}$  where  $\mathbf{d}$  is any linear combination of the vectors  $\mathbf{b}_J, \mathbf{b}_{J+1}, \dots, \mathbf{b}_p$ . A similar result will hold for any set of basis vectors of the vector space generated by the columns of the matrix

$$\mathbf{B}^* = [\mathbf{b}_J, \mathbf{b}_{J+1}, \dots, \mathbf{b}_p], \tag{21}$$

so that  $\bar{\mathbf{X}}\mathbf{B}^* = \mathbf{0}$ .

Therefore, a set of  $r - J + 1$  linear independent vectors of the form  $\mathbf{d}$  where  $\mathbf{d}$  is a linear combination of any basis of  $\mathcal{V}(\mathbf{B}^*)$  is needed such that the scaffolding vectors consisting of the first  $J - 1$  columns of  $\mathbf{B}$  together with the  $r - J + 1$   $\mathbf{d}$  vectors minimize  $TSRES$  for all legitimate choices of the  $\{\mathbf{d}\}$ . Write these  $r$  scaffolding vectors as the columns of the matrix  $\mathbf{D}_r$  i.e.,

$$\mathbf{D}_r = [b_1, b_2, \dots, b_{J-1}, \mathbf{d}_J, \mathbf{d}_{J+1}, \dots, \mathbf{d}_r] \tag{22}$$

and let the columns of  $\mathbf{D} : p \times p = [\mathbf{D}_r, \mathbf{d}_{r+1}, \dots, \mathbf{d}_p]$  represent a basis of  $\mathcal{V}(\mathbf{B})$ . It follows that any column of  $\mathbf{B}$  can be written as a linear combination of the columns

of  $D$ . Therefore, there exists a non-singular matrix  $C : p \times p$  such that  $B = DC$  i.e.,  $D = BF$  with  $F = C^{-1}$ .

Straightforward algebraic manipulation shows that  $F$  is of the form

$$F = \begin{bmatrix} I_{J-1} & \mathbf{0} \\ \mathbf{0} & F^* \end{bmatrix}, \tag{23}$$

where  $F^*$  is an  $(p - J + 1) \times (p - J + 1)$  orthogonal matrix. We provide a detailed derivation as supplementary material. Write

$$F^* = [f_1^*, f_2^*, \dots, f_{p-J+1}^*]. \tag{24}$$

Then,  $F^{-1} = \begin{bmatrix} I & \mathbf{0} \\ \mathbf{0} & (F^*)^T \end{bmatrix}$  and our scaffolding vectors for constructing the  $r$ -dimensional CVA biplot are the columns of

$$\begin{aligned} D_r &= [b_1, b_2, \dots, b_{J-1}, d_J, d_{J+1}, \dots, d_r] \\ &= B \left[ \begin{bmatrix} I_{J-1} \\ \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ f_1^* \end{bmatrix} \cdots \begin{bmatrix} \mathbf{0} \\ f_{r-J+1}^* \end{bmatrix} \right]. \end{aligned} \tag{25}$$

Furthermore,

$$D^{(r)} = \begin{bmatrix} I_{J-1} & \mathbf{0} \\ \mathbf{0}^T & (f_1^*)^T \\ \dots & \dots \\ \mathbf{0}^T & (f_{r-J+1}^*)^T \end{bmatrix} B^{-1}. \tag{26}$$

The approximations of the rows of  $X$ , i.e., the original samples, in the biplot constructed on the scaffolding provided by the columns of  $D_r$  follow from (16) and using (25) and (26) as

$$\begin{aligned} \hat{X} &= X D_r D^{(r)} \\ &= X B \left[ \begin{bmatrix} I_{J-1} \\ \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{0} \\ f_1^* \end{bmatrix} \cdots \begin{bmatrix} \mathbf{0} \\ f_{r-J+1}^* \end{bmatrix} \right] \begin{bmatrix} I_{J-1} & \mathbf{0} \\ \mathbf{0}^T & (f_1^*)^T \\ \dots & \dots \\ \mathbf{0}^T & (f_{r-J+1}^*)^T \end{bmatrix} B^{-1} \\ &= X B \begin{bmatrix} I_{J-1} & \mathbf{0} \\ \mathbf{0}^T & f_1^*(f_1^*)^T + f_2^*(f_2^*)^T + \dots + f_{r-J+1}^*(f_{r-J+1}^*)^T \end{bmatrix} B^{-1}, \end{aligned} \tag{27}$$

where, from the orthogonality of  $F^*$ , it follows that  $(f_1^*)^T f_1^* = (f_2^*)^T f_2^* = \dots = (f_{r-J+1}^*)^T f_{r-J+1}^* = 1$ .

The criterion *TSRES* now becomes

$$\begin{aligned} & \|X - \hat{X}\|^2 \\ &= tr\{(X - \hat{X})(X - \hat{X})^T\} \\ &= \left\| X \left( I_p - B \begin{bmatrix} I_{J-1} & \mathbf{0} \\ \mathbf{0}^T & f_1^*(f_1^*)^T + f_2^*(f_2^*)^T + \dots + f_{r-J+1}^*(f_{r-J+1}^*)^T \end{bmatrix} B^{-1} \right) \right\|^2. \end{aligned} \tag{28}$$

To construct an  $r$ -dimensional CVA biplot satisfying our aim of minimizing *TSRES* while simultaneously providing 100% accurate predictions for the  $J$  sample means when  $J < p$ , we propose the following:

- Find the solution of

$$argmin \left\| X \left( I_p - BLB^{-1} \right) \right\|^2, \tag{29}$$

where  $L = \begin{bmatrix} I_{J-1} & \mathbf{0} \\ \mathbf{0}^T & f_1^*(f_1^*)^T + f_2^*(f_2^*)^T + \dots + f_{r-J+1}^*(f_{r-J+1}^*)^T \end{bmatrix}$  and the minimum is taken with respect to the  $f_j^*$  such that  $(f_j^*)^T f_j^* = 1$  for  $j = 1, 2, \dots, r - J + 1$ .

- Use the optimum  $\{f_1^*, f_2^*, \dots, f_{r-J+1}^*\}$  to construct  $(B_{opt})_r = [b_1, \dots, b_{J-1}, d_J, \dots, d_r]$  where

$$\begin{aligned} d_j &= B^*(f_{j-J+1}^*)_{opt} \\ &= B \begin{bmatrix} \mathbf{0} \\ (f_{j-J+1}^*)_{opt} \end{bmatrix} \\ &= f_{J(j-J+1)}^{opt} b_J + f_{(J+1)(j-J+1)}^{opt} b_{J+1} + \dots + f_{p(j-J+1)}^{opt} b_p, \end{aligned}$$

for  $j = J, J + 1, \dots, r$ .

- Next,  $(B_{opt})_r$  is used for constructing the  $r$ -dimensional CVA biplot with calibrated prediction (or interpolation) axes.
- Finally, calculate a standardised form of  $min(TSRES)$ :  $\frac{min(TSRES(X, \hat{X}))}{tr(XX^T)}$ , as a measure of the accuracy of the approximations of the individual sample points in the  $r$ -dimensional biplot.

The solution for (29) can be found from (28) as follows: Since  $B^T S_W B = I_p$  it follows that

$$\begin{aligned} I_p &= B^T X^T X B - B^T \bar{X}^T N_g \bar{X} B \\ &= \begin{bmatrix} (B_{J-1})^T X^T X B_{J-1} & (B_{J-1})^T X^T X B^* \\ (B^*)^T X^T X B_{J-1} & (B^*)^T X^T X B^* \end{bmatrix} \end{aligned}$$



$$\begin{aligned}
 & - \begin{bmatrix} (\mathbf{B}_{J-1})^T \bar{\mathbf{X}}^T \mathbf{N}_g \bar{\mathbf{X}} \mathbf{B}_{J-1} & (\mathbf{B}_{J-1})^T \bar{\mathbf{X}}^T \mathbf{N}_g \bar{\mathbf{X}} \mathbf{B}^* \\ (\mathbf{B}^*)^T \bar{\mathbf{X}}^T \mathbf{N}_g \bar{\mathbf{X}} \mathbf{B}_{J-1} & (\mathbf{B}^*)^T \bar{\mathbf{X}}^T \mathbf{N}_g \bar{\mathbf{X}} \mathbf{B}^* \end{bmatrix} \\
 = & \begin{bmatrix} (\mathbf{B}_{J-1})^T \mathbf{X}^T \mathbf{X} \mathbf{B}_{J-1} & (\mathbf{B}_{J-1})^T \mathbf{X}^T \mathbf{X} \mathbf{B}^* \\ (\mathbf{B}^*)^T \mathbf{X}^T \mathbf{X} \mathbf{B}_{J-1} & (\mathbf{B}^*)^T \mathbf{X}^T \mathbf{X} \mathbf{B}^* \end{bmatrix} - \begin{bmatrix} (\mathbf{B}_{J-1})^T \bar{\mathbf{X}}^T \mathbf{N}_g \bar{\mathbf{X}} \mathbf{B}_{J-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix},
 \end{aligned}$$

establishing that

$$(\mathbf{B}^*)^T \mathbf{X}^T \mathbf{X} \mathbf{B}^* = \mathbf{I}_{p-J+1}. \tag{30}$$

Since  $\mathbf{B} = [\mathbf{B}_{J-1} \ \mathbf{B}^*]$  and non-singular, we set

$$\mathbf{B}^{-1} = \begin{bmatrix} \mathbf{B}^{(J-1)} : (J-1) \times p \\ \mathbf{B}^{(2)} : (p-J+1) \times p \end{bmatrix}.$$

Write  $\mathbf{U} : (p-J+1) \times (p-J+1) = \mathbf{f}_1^*(\mathbf{f}_1^*)^T + \mathbf{f}_2^*(\mathbf{f}_2^*)^T + \dots + \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T$ ; then it follows that

$$\begin{aligned}
 & \mathbf{X} \left( \mathbf{I}_p - \mathbf{B} \begin{bmatrix} \mathbf{I}_{J-1} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{f}_1^*(\mathbf{f}_1^*)^T + \mathbf{f}_2^*(\mathbf{f}_2^*)^T + \dots + \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \end{bmatrix} \mathbf{B}^{-1} \right) \\
 & = \mathbf{X} \left( \mathbf{I}_p - \left( \mathbf{B}_{J-1} \mathbf{B}^{(J-1)} + \mathbf{B}^* \mathbf{U} \mathbf{B}^{(2)} \right) \right). \tag{31}
 \end{aligned}$$

From  $[\mathbf{B}_{J-1} \ \mathbf{B}^*] \begin{bmatrix} \mathbf{B}^{(J-1)} \\ \mathbf{B}^{(2)} \end{bmatrix} = \mathbf{I}_p$  it follows that  $\mathbf{B}_{J-1} \mathbf{B}^{(J-1)} = \mathbf{I}_p - \mathbf{B}^* \mathbf{B}^{(2)}$  so that (31) can be written as

$$\mathbf{X} \left( \mathbf{I}_p - \left( \mathbf{B}_{J-1} \mathbf{B}^{(J-1)} + \mathbf{B}^* \mathbf{U} \mathbf{B}^{(2)} \right) \right) = \mathbf{X} \left( \mathbf{B}^* \mathbf{B}^{(2)} - \mathbf{B}^* \mathbf{U} \mathbf{B}^{(2)} \right). \tag{32}$$

Therefore,

$$\begin{aligned}
 & \left\| \mathbf{X} \left( \mathbf{I}_p - \mathbf{B} \begin{bmatrix} \mathbf{I}_{J-1} & \mathbf{0} \\ \mathbf{0}^T & \mathbf{f}_1^*(\mathbf{f}_1^*)^T + \mathbf{f}_2^*(\mathbf{f}_2^*)^T + \dots + \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \end{bmatrix} \mathbf{B}^{-1} \right) \right\|^2 \\
 & = \left\| \mathbf{X} \mathbf{B}^* (\mathbf{I}_{p-J+1} - \mathbf{U}) \mathbf{B}^{(2)} \right\|^2 \\
 & = \text{tr} \left\{ \mathbf{X} \mathbf{B}^* (\mathbf{I}_{p-J+1} - \mathbf{U}) \mathbf{B}^{(2)} (\mathbf{B}^{(2)})^T (\mathbf{I}_{p-J+1} - \mathbf{U}) (\mathbf{B}^*)^T \mathbf{X}^T \right\} \\
 & = \text{tr} \left\{ \mathbf{X} \mathbf{B}^* (\mathbf{B}^{(2)})^T (\mathbf{B}^*)^T \mathbf{X}^T \right\} - 2 \text{tr} \left\{ \mathbf{X} \mathbf{B}^* \mathbf{U} \mathbf{B}^{(2)} (\mathbf{B}^{(2)})^T (\mathbf{B}^*)^T \mathbf{X}^T \right\} \\
 & \quad + \text{tr} \left\{ \mathbf{X} \mathbf{B}^* \mathbf{U} (\mathbf{B}^{(2)})^T \mathbf{U} (\mathbf{B}^*)^T \mathbf{X}^T \right\} \tag{33}
 \end{aligned}$$

$$= \text{tr} \left\{ \mathbf{B}^{(2)} (\mathbf{B}^{(2)})^T \right\} - 2 \text{tr} \left\{ \mathbf{U} \mathbf{B}^{(2)} (\mathbf{B}^{(2)})^T \right\} + \text{tr} \left\{ \mathbf{U} \mathbf{B}^{(2)} (\mathbf{B}^{(2)})^T \mathbf{U} \right\} \tag{34}$$

by substituting (30) in (33). Write  $\mathbf{H} = \mathbf{B}^{(2)}(\mathbf{B}^{(2)})^T : (p - J + 1) \times (p - J + 1)$  then it follows that  $\mathbf{H}$  has rank  $(p - J + 1)$  and is thus positive definite. Therefore,

$$\begin{aligned} & \left\| \mathbf{X} \left( \mathbf{I}_p - \mathbf{B} \begin{bmatrix} \mathbf{I}_{J-1} & & & \mathbf{0} \\ \mathbf{0}^T & \mathbf{f}_1^*(\mathbf{f}_1^*)^T & & \\ & \mathbf{f}_2^*(\mathbf{f}_2^*)^T & & \\ & & \ddots & \\ & & & \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \end{bmatrix} \mathbf{B}^{-1} \right) \right\|^2 \\ & = \text{tr}\{\mathbf{H}\} - 2\text{tr}\{\mathbf{UH}\} + \text{tr}\{\mathbf{UHU}\}, \end{aligned} \tag{35}$$

where

$$\text{tr}\{\mathbf{UH}\} = (\mathbf{f}_1^*)^T \mathbf{H} \mathbf{f}_1^* + (\mathbf{f}_2^*)^T \mathbf{H} \mathbf{f}_2^* + \dots + (\mathbf{f}_{r-J+1}^*)^T \mathbf{H} \mathbf{f}_{r-J+1}^* \tag{36}$$

and

$$\begin{aligned} \text{tr}\{\mathbf{UHU}\} &= \text{tr}\{\mathbf{UUH}\} \\ &= \text{tr}\left\{ \left( \mathbf{f}_1^*(\mathbf{f}_1^*)^T + \mathbf{f}_2^*(\mathbf{f}_2^*)^T + \dots + \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \right) \right. \\ & \quad \left. \times \left( \mathbf{f}_1^*(\mathbf{f}_1^*)^T + \mathbf{f}_2^*(\mathbf{f}_2^*)^T + \dots + \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \right) \mathbf{H} \right\}. \end{aligned}$$

Since  $(\mathbf{f}_i^*)^T \mathbf{f}_j^* = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$  it follows that

$$\begin{aligned} \text{tr}\{\mathbf{UHU}\} &= \text{tr}\left\{ \left( \mathbf{f}_1^*(\mathbf{f}_1^*)^T \mathbf{H} + \mathbf{f}_2^*(\mathbf{f}_2^*)^T \mathbf{H} + \dots + \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \mathbf{H} \right) \right\} \\ &= (\mathbf{f}_1^*)^T \mathbf{H} \mathbf{f}_1^* + (\mathbf{f}_2^*)^T \mathbf{H} \mathbf{f}_2^* + \dots + (\mathbf{f}_{r-J+1}^*)^T \mathbf{H} \mathbf{f}_{r-J+1}^*. \end{aligned} \tag{37}$$

Substituting (36) and (37) into (35) leads to

$$\begin{aligned} & \left\| \mathbf{X} \left( \mathbf{I}_p - \mathbf{B} \begin{bmatrix} \mathbf{I}_{J-1} & & & \mathbf{0} \\ \mathbf{0}^T & \mathbf{f}_1^*(\mathbf{f}_1^*)^T & & \\ & \mathbf{f}_2^*(\mathbf{f}_2^*)^T & & \\ & & \ddots & \\ & & & \mathbf{f}_{r-J+1}^*(\mathbf{f}_{r-J+1}^*)^T \end{bmatrix} \mathbf{B}^{-1} \right) \right\|^2 \\ & = \text{tr}\{\mathbf{H}\} - \left( (\mathbf{f}_1^*)^T \mathbf{H} \mathbf{f}_1^* + (\mathbf{f}_2^*)^T \mathbf{H} \mathbf{f}_2^* + \dots + (\mathbf{f}_{r-J+1}^*)^T \mathbf{H} \mathbf{f}_{r-J+1}^* \right). \end{aligned} \tag{38}$$

Remembering that  $\mathbf{H}$  is positive definite, the criterion (38) can be minimized by maximizing each of the terms  $(\mathbf{f}_j^*)^T \mathbf{H} \mathbf{f}_j^*$  with respect to  $\mathbf{f}_j^*$  for  $j = 1, 2, \dots, r - J + 1$  under the constraint  $(\mathbf{f}_j^*)^T \mathbf{f}_j^* = 1$ . This is readily accomplished by introducing the Lagrange multiplier  $\lambda_j$  to form

$$(\mathbf{f}_j^*)^T \mathbf{H} \mathbf{f}_j^* - \lambda_j \left( (\mathbf{f}_j^*)^T \mathbf{f}_j^* - 1 \right). \tag{39}$$

Differentiating (39) with respect to  $\lambda_j$  and equating to zero leads to  $\mathbf{H} \mathbf{f}_j^* = \lambda_j \mathbf{f}_j^*$  i.e.,  $\lambda_j = (\mathbf{f}_j^*)^T \mathbf{H} \mathbf{f}_j^*$ . Thus  $(\mathbf{f}_j^*)^T \mathbf{H} \mathbf{f}_j^*$  is maximized when  $\mathbf{f}_j^*$  is a normalized eigenvector associated with the  $j$ th largest eigenvalue of  $\mathbf{H}$  and hence the criterion (39) attains its minimum when the  $\{\mathbf{f}_j^*\}$  are set to the normalized eigenvectors associated with the largest  $r - J + 1$  eigenvalues of  $\mathbf{H} = \mathbf{B}^{(2)}(\mathbf{B}^{(2)})^T$ , respectively. Denote

these eigenvectors by  $f_1^{opt}, f_2^{opt}, \dots, f_{r-J+1}^{opt}$ , respectively, where  $r \leq p$  then a  $p \times r$  matrix  $(\mathbf{B}_{opt})_r$  can be constructed as

$$(\mathbf{B}_{opt})_r = \left[ \mathbf{B}_{J-1}, \mathbf{B}^* \left[ f_1^{opt}, f_2^{opt}, \dots, f_{r-J+1}^{opt} \right] \right]. \tag{40}$$

Setting  $r = p$  in the above leads to a matrix  $\mathbf{B}_{opt}$  of size  $p \times p$  which is non-singular, allowing for the computation of the matrices  $(\mathbf{B}_{opt})_r$ , consisting of the first  $r$  columns of  $\mathbf{B}_{opt}$ , and  $(\mathbf{B}_{opt})^{(r)}$ , consisting of the first  $r$  rows of  $(\mathbf{B}_{opt})^{-1}$ . Therefore

$$\hat{\mathbf{X}} = \mathbf{X} \mathbf{D}_r \mathbf{D}^{(r)} = \mathbf{X} (\mathbf{B}_{opt})_r (\mathbf{B}_{opt})^{(r)} \tag{41}$$

will minimize *TSRES*.

### 6 An alternative biplot based on the Bhattacharyya distance for the two sample case

If, analogous to section 3 the population parameters in (5) are replaced with their sample estimates the sample version of the Bhattacharyya distance consists of two terms. The first measures the dissimilarity between the two sample means and the second measures the dissimilarity between the two sample covariance matrices. Fukunaga (1990) uses this property to construct a second dimension for a visual display of the rows of a data matrix  $\mathbf{X} : n \times p$  in two dimensions. Furthermore, Hennig (2004) utilizes the Bhattacharyya distance, among other methods to construct two-dimensional visualizations of the dispersions of two asymmetric samples – asymmetric in the sense that one is known to be more homogeneous and the other to be more heterogeneous. However, it should be noted that none of the visualizations proposed by Fukunaga (1990) and Hennig (2004) are biplots because no information regarding the columns of  $\mathbf{X}$  is displayed. The optimal two-group CVA biplot proposed in section 5 assumes equality of covariance matrices as is usual for CVA. This implies that the second term of the sample version of (5) will vanish and optimization of (5) becomes equivalent to maximizing (9).

Denote the transformation to the canonical space  $\mathcal{C}$  by  $\mathbf{Y} = \mathbf{X} \mathbf{M}$  with  $\mathbf{M} : p \times p$ . We can write

$$\mathcal{C} = \mathcal{C}_1 \cup \mathcal{V},$$

where  $\mathcal{C}_1$  is one-dimensional based on  $\mathbf{m}_1$  and  $\mathcal{V}$  is  $(p - 1)$ -dimensional with basis  $\mathbf{m}_2, \mathbf{m}_3, \dots, \mathbf{m}_p$ . In  $\mathcal{V}$  we have

$$\begin{aligned} \mathbf{Y}^{*T} &= [Y_2, Y_3, \dots, Y_p], \\ \mathbf{Y}^* | G_i : (p - 1) \times 1 &\sim (\boldsymbol{\mu}^*, \boldsymbol{\Sigma}_i^*), i = 1; 2, \end{aligned}$$

with  $\boldsymbol{\mu}^{*T} = [\mu_2^{(Y)}, \mu_3^{(Y)}, \dots, \mu_p^{(Y)}]$ ,  $\boldsymbol{\Sigma}_1^* = \mathbf{M}^{*T} \boldsymbol{\Sigma}_1 \mathbf{M}^*$ ,  $\boldsymbol{\Sigma}_2^* = \mathbf{M}^{*T} \boldsymbol{\Sigma}_2 \mathbf{M}^*$  and  $\mathbf{M}^* : p \times (p - 1) = [\mathbf{m}_2, \mathbf{m}_3, \dots, \mathbf{m}_p]$ .

Since  $\mu_1^* = \mu_2^* = \mu^*$ , (say) we have from (5)

$$D_{Bhat} = \frac{1}{2} \log \left( \frac{\left| \frac{\Sigma_1^* + \Sigma_2^*}{2} \right|}{\sqrt{|\Sigma_1^*| |\Sigma_2^*|}} \right). \tag{42}$$

Maximization of the sample version of (42) proceeds parallel to the process described in section 2.3. The outcome is the matrix  $A : (p - 1) \times (p - 1)$  containing as columns the required eigenvectors arranged in decreasing order of the values of  $\lambda_i^* + \frac{1}{\lambda_i^*} + 2$ .

A 2D biplot can now be constructed using the methods described in sections 3 and 4 by first noting that the matrix  $M$  is available as the matrix  $B$  of section 3. Next, we calculate the matrix

$$K : p \times p = [m_1 \quad M^*A]$$

and its inverse  $K^{-1}$ . Let  $a : (p - 1) \times 1$  denote the first column of  $A$ , then the 2D biplot can be constructed as described in section 4 using the rows of  $Z : n \times 2$ , where

$$Z = [Xm_1 \quad Y^*a] = X[m_1 \quad M^*a]$$

for plotting the samples. The variables are represented by  $p$  calibrated biplot axes constructed from

$$\frac{\text{marker}}{e_k^T K^{(2)T} K^{(2)} e_k} K^{(2)} e_k,$$

where  $K^{(2)}$  denotes the first 2 rows of  $K^{-1}$ .

## 7 Examples

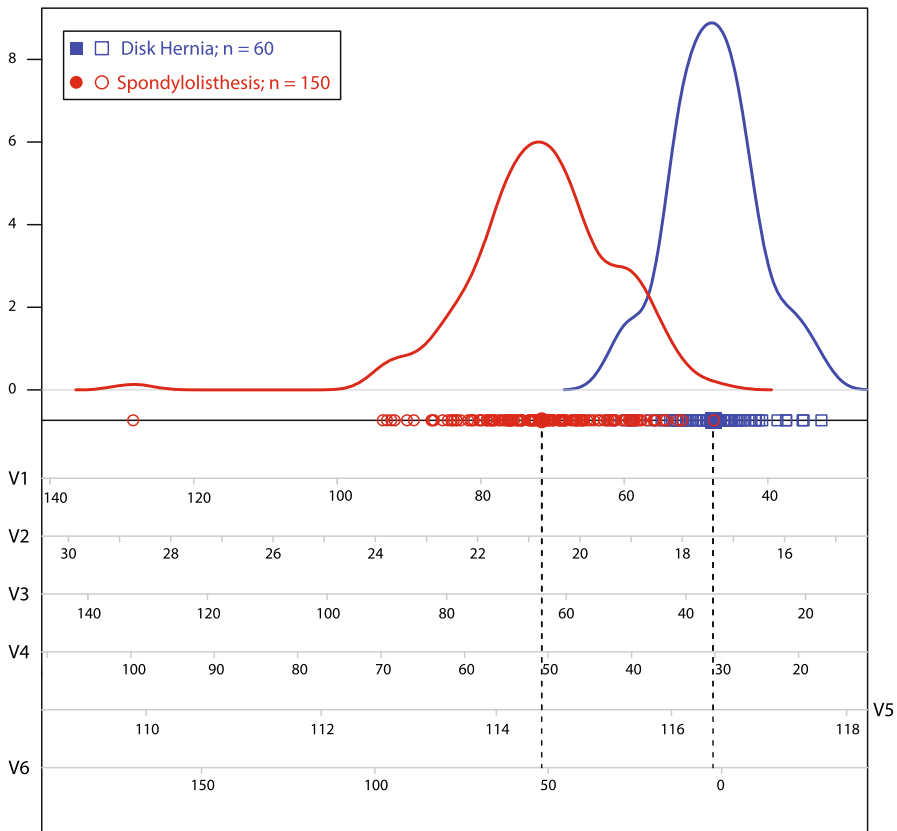
As an example, we consider the Vertebral Column Data Set from the UCI Machine Learning Repository (Barreto et al. 2011) and discussed in detail by da Rocha Neto et al. (2011). The full data set contains measurements on six continuous/numeric variables relating to the shape and orientation of the pelvis and lumbar spine for each of the 310 individuals (samples). These samples are classified as normal, disk hernia, or spondylolisthesis patients. As an example of a two-group CVA, we study the subset of 60 disk hernia and 150 spondylolisthesis patients. The six numeric variables are Pelvic incidence, Pelvic tilt, Lumbar lordosis angle, Sacral slope, Pelvic radius, and Degree spondylolisthesis.

Table 1 contains the group means for each of the six variables - with the outlier (see Fig. 2) included and excluded. It is hard to see from the table the presence of the outlier but, as is evident from Figs. 1 and 2, the outlier will be clearly revealed in a CVA biplot.

The one-dimensional biplot of this two-group data set is shown in Fig. 1. In this biplot, the two group means as well as all six variables (in the form of six calibrated axes) lie on the one scaffolding line defined by the singular vector associated with

**Table 1** Group means for the six measurements V1 = Pelvic incidence, V2 = Pelvic tilt, V3 = Lumbar lordosis angle, V4 = Sacral slope, V5 = Pelvic radius and V6 = Degree spondylolisthesis

	Disk Hernia	Spondylolisthesis	
		Outlier included	Outlier excluded
V1	47.6383	71.5137	71.1223
V2	17.3987	20.7480	20.8309
V3	35.4635	64.1099	64.2154
V4	30.2400	50.7661	50.2919
V5	116.4750	114.5183	114.5642
V6	2.4793	51.8969	49.4362



**Fig. 1** 1D biplot of the two-group Vertebral Column data. V1 = Pelvic incidence, V2 = Pelvic tilt, V3 = Lumbar lordosis angle, V4 = Sacral slope, V5 = Pelvic radius, V6 = Degree spondylolisthesis. The one-dimensional biplot is enhanced by superimposing density estimates for the interpolated sample points in the two groups

the single non-zero singular value of the underlying two-sided eigenequation. The six calibrated axes representing the variables have been vertically translated to aid the interpretation of the biplot. The values for each of the group means for all six variables can easily be read from the biplot axes and it can be verified that these values coincide exactly with the corresponding values in Table 1. As expected all variables have predictivities of 100% for determining the mean values with  $TSREM = 0$ . All the individual sample points have been interpolated onto the biplot and therefore they also lie on the single scaffolding axis defining the biplot. However,  $TSRES > 0$  with a standardized version  $TSRES/(tr(\mathbf{X}\mathbf{X}^T)) = 0.4702$ . To visualize the within groups variation the biplot has been enhanced by the addition of density estimates of the two sets of sample points interpolated onto the one-dimensional CVA biplot. These density estimates show graphically the separation/overlap of the two groups. Inspection of the six biplot axes suggests V5 to be negatively correlated with the other five variables; the latter are all pairwise positively correlated.

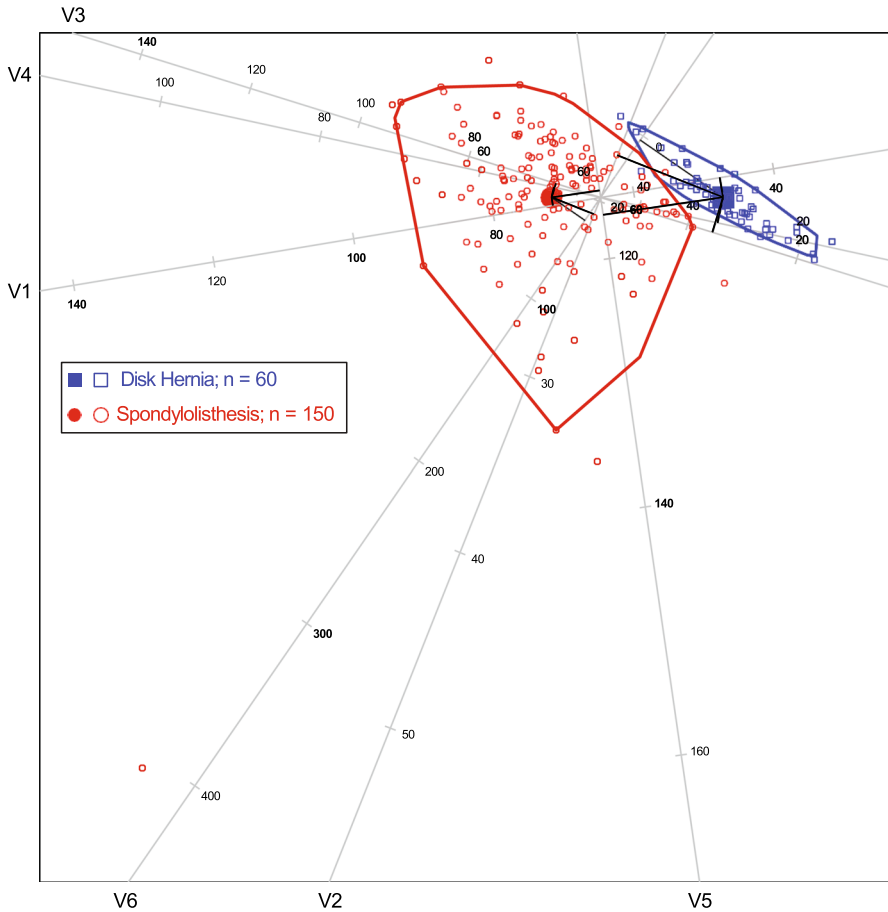
It is clear that much can be learned from the one-dimensional CVA biplot but some serious issues are calling for considering a second scaffolding axis:

- Is the conspicuous outlier an outlier on all variables?
- Can the approximation of the sample points be improved without sacrificing what has been achieved with the mean vectors?
- Is it possible to construct a more detailed visualization of the separation/overlap of the two groups?
- Is it possible to construct a more accurate and detailed visualization of the correlation structure of the six variables?

Figure 2 provides an answer to the above issues. The optimal two-dimensional CVA biplot in Fig. 2 demonstrates the following:

- The biplot is uniquely defined.
- Each of the six biplot axes has a predictivity of 100% for determining the values of the two groups for all variables. This results in  $TSREM$  to remain zero. Figure 2 provides a clearer picture of how the predictivities are determined.
- Introducing the second scaffolding axis improves the approximation of the samples appreciably: the standardized  $TSRES$  decreases to 0.1799 (a decrease of more than 60% of the corresponding value obtained in Fig. 2).
- The second scaffolding axis is optimal in the sense that no other scaffolding axis can be found, which will improve  $TSRES$  while restricting  $TSREM$  to zero.
- It is clear that Sample 116 is less of an outlier on V3 and V4 than on V2 and V6.
- There is a suggestion that while V5 appears to be negatively correlated with V3 and V4 it appears to be positively correlated with V2 and V6. We note that the addition of a second scaffolding axis provides angles between the biplot axes, which allow for visualizing the approximate correlational structure between the variables.

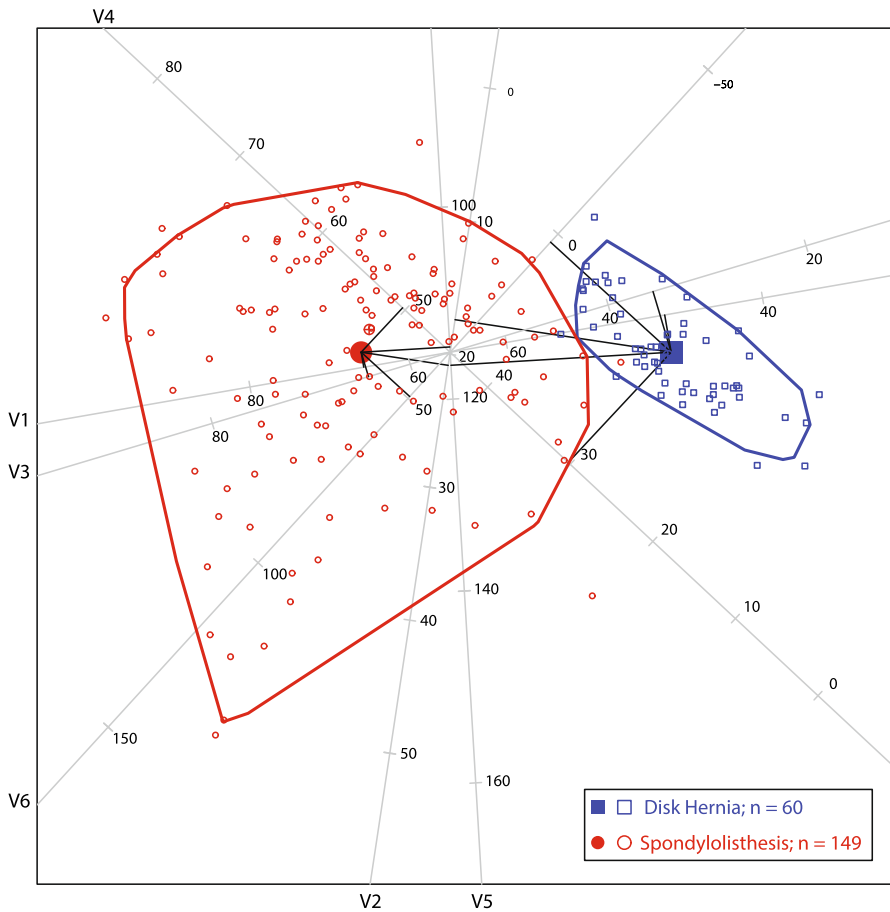
The biplot in Fig. 2 has been enhanced by superimposing 95%-bags onto the biplot. Alpha-bags are discussed in detail by Gower et al. (2011). The 95%-bag used here contains the innermost 95% of the bivariate sample points, where the innermost is relative to the Tukey median (Ruts and Rousseeuw 1996). Now, we are ready for a detailed appraisal of the overlap/separation of the two groups based on the two-dimensional



**Fig. 2** Optimal 2D CVA biplot of the two-group Vertebral Column data. V1 = Pelvic incidence, V2 = Pelvic tilt, V3 = Lumbar lordosis angle, V4 = Sacral slope, V5 = Pelvic radius, V6 = Degree spondylolisthesis. The optimal two-dimensional biplot is enhanced by superimposing bags containing the innermost 95% samples of the two groups respectively

clouds of points visualizing the within groups sample variation. However, we first exclude Sample 116 from the analysis and consider the optimal two-dimensional CVA biplot given in Fig. 3 overlaid with 95%-bags. This figure shows:

- Clearly how the biplot axes are used to determine the group means exactly for each variable.
- The angles between the biplot axes allow for an approximate visual appraisal of the correlation structure.
- It is seen that the two 95%-bags almost touch each other but do not overlap giving us an overall quantitative measure of the degree of separation between the two groups.
- The standardized *TSRES* value is 0.2155.
- Although the two clouds of points have a high degree of separation it can also be seen that



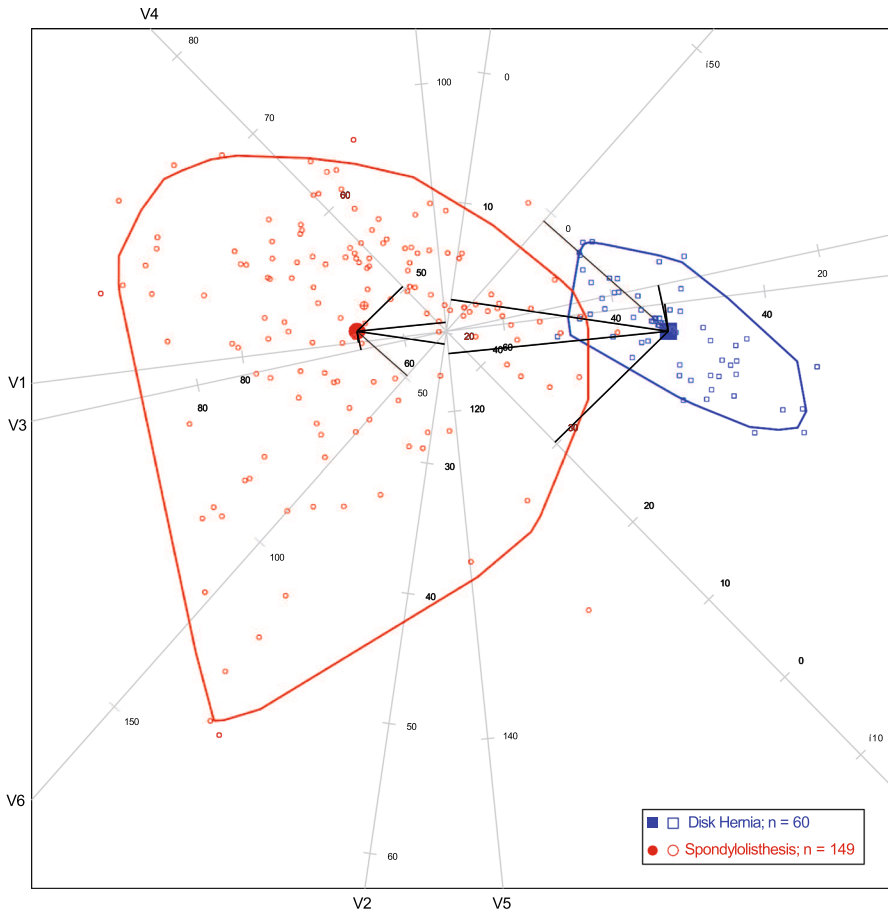
**Fig. 3** Optimal 2D CVA biplot of the two-group Vertebral Column data excluding outlier sample 116. V1 = Pelvic incidence, V2 = Pelvic tilt, V3 = Lumbar lordosis angle, V4 = Sacral slope, V5 = Pelvic radius, V6 = Degree spondylolisthesis. The optimal two-dimensional biplot is enhanced by superimposing bags containing the innermost 95% samples of the two groups respectively

- there is a high degree of overlap between the two groups concerning V2 and V5;
- there is almost no overlap on V1, V3 and V6;
- on V4 large Disk Hernia values overlap with small Spondylolisthesis measurements, while small measurements of V4 almost exclusively occur in Disk Hernia.

The CVA biplots in Figs. 2 and 3 assume equal covariance matrices. We can now relax this assumption and construct in Fig. 4, the biplot based on the Bhattacharyya distance as discussed in section 6.

Although the appearance of this biplot is quite similar to that of the corresponding optimal CVA biplot shown in Fig. 3 its standardized *TSRES* value is approximately 10% higher, namely 0.2367.





**Fig. 4** The 2D biplot based on the Bhattacharyya distance of the two-group Vertebral Column data excluding outlier sample 116. V1 = Pelvic incidence, V2 = Pelvic tilt, V3 = Lumbar lordosis angle, V4 = Sacral slope, V5 = Pelvic radius, V6 = Degree spondylolisthesis. The optimal two-dimensional biplot is enhanced by superimposing bags containing the innermost 95% samples of the two groups respectively

## 8 Conclusions

CVA biplots are constructed to show in a single plot the group means as points and all the variables as calibrated linear biplot axes. In the case of two groups the CVA biplot becomes a line containing all these points and biplot axes. Furthermore, there is no approximation in the positions of the group means and their respective values, which can be exactly determined from the biplot axes. It is common practice to interpolate the individual sample points onto the CVA biplot as well, but then they appear as approximations in the one-dimensional CVA biplot space. Since all the biplot axes lie on top of each other it is difficult to use them for determining the values of the two means for the different variables. However, the vertical translation of one of these axes

does not affect the values of the two means for that particular variable. Therefore, as has been shown in Fig. 1, vertical translation of the biplot axes does not change the dimensionality of the CVA biplot but increases the usefulness of the different biplot axes for reading off values of the respective variables.

The fundamental question that is addressed in this paper is: What can be gained, if anything, by increasing the dimensionality of the above one-dimensional CVA biplot to two dimensions? This question can be rephrased as: Can we add a second dimension to our one-dimensional CVA biplot to improve the approximation of the individual sample points leading to a better understanding of the within groups variability while leaving unchanged the optimal representation of the two group means? As it turned out the addition of a second dimension is not a straightforward process since there are an infinite number of ways that this can be done. Furthermore, if existing software is used for constructing a two-dimensional CVA biplot in the two-group case the result can be highly misleading. This is because the two-sided eigenequation underlying the CVA procedure has only a single non-zero singular value with no natural ordering of the zero singular values resulting in the indeterminacy of singular vectors associated with zero singular values. Therefore, to guarantee a unique solution for finding a second dimension, we had to consider a criterion to optimize the approximation of the sample points while leaving the optimal representation of the group means unchanged. We suggested the *TSRES* criterion for meeting this goal. Minimizing *TSRES* results in a uniquely defined two-dimensional CVA biplot for the two-group case. It optimizes the approximations of the within-group samples while the two group means are exactly represented with the linear biplot axes providing exact values for both groups on all variables.

Figure 3 shows that our proposed optimal two-dimensional CVA biplot for two groups has the potential to provide the researcher with a tool that not only distinguishes the group means optimally but also where the within sample variation can be depicted graphically to gain deeper insight into the separation/overlap of the two groups. Moreover, it is unique and thus prevents any possibility of ambiguity when routinely using existing software. Thus we have attained our primary objective as is illustrated in the example discussed above.

The algebra underlying the optimal two-dimensional CVA biplot extends directly to an optimal three-dimensional biplot when dealing with a three-groups case having a CVA biplot space of dimension two.

An alternative suggestion based on the Bhattacharyya distance is available when relaxing the equal within-group covariance matrix assumption. As can be seen from Fig. 4, the biplot is slightly different, but the overall conclusion regarding overlap and separation in terms of the individual variables remains unchanged. However, the biplot based on Bhattacharyya distance is designed to optimize a different objective than the optimization of the sum of the squared approximations of the data matrix. Therefore, our preferred 2D CVA biplot to construct in the two-group case is the proposed optimal CVA biplot.

Finally, we note that our R code for constructing the biplots discussed in this paper is available in Lubbe et al. (2023).

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1007/s11634-024-00593-7>.

**Acknowledgements** We are grateful to the comments made by an anonymous reviewer and the review editor that have improved the quality of this manuscript.

**Funding** Open access funding provided by Stellenbosch University. This work is based upon research supported in part by the National Research Foundation (NRF) of South Africa (Grant Number 103310). Any opinions, findings and conclusions, or recommendations expressed in this material are those of the authors, and therefore the NRF does not accept any liability in regard thereof.

## Declarations

**Conflict of interest** The authors declare that they have no Conflict of interest.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

- Barreto GdA, da Rocha Neto AR, da Mota Filho HAF (2011) Vertebral column data set. UCI Machine Learning Repository, <http://archive.ics.uci.edu/ml>
- Fisher RA (1936) The use of multiple measurements in taxonomic problems. *Annals of Eugenics* 7:179–188
- Flury B (1997) A first course in multivariate statistics. Springer-Verlag, New York
- Flury L, Boukai B, Flury BD (1997) The discrimination subspace model. *J American Stat Assoc* 92:758–766
- Fukunaga K (1990) Introduction to statistical pattern recognition, 2nd edn. Academic Press, Boston
- Gabriel KR (1971) The biplot graphical display of matrices with application to principal component analysis. *Biometrika* 58:451–467
- Gabriel KR (1972) Analysis of meteorological data by means of manonical decomposition and biplots. *J Appl Meteorol* 11:1071–1077
- Gardner-Lubbe S, Le Roux NJ, Gower JC (2008) Measures of fit in principal component and canonical variate analyses. *J Appl Stat* 35:947–965
- Gittins R (1985) Canonical analysis. Springer-Verlag, Berlin
- Gower JC (1995) A general theory of biplots. In: Krzanowski WJ (ed) Recent advances in descriptive multivariate analysis. Clarendon Press, Oxford, pp 283–303
- Gower JC, Hand DJ (1996) Biplots. Chapman & Hall, London
- Gower JC, Le Roux NJ, Lubbe S (2011) Understanding biplots. Wiley, Chichester
- Hastie T, Tibshirani R, Friedman J (2001) The elements of statistical learning. Springer-Verlag, New York
- Hennig C (2004) Asymmetric linear dimension reduction for classification. *J Comput Graph Stat* 13:930–945
- Lubbe S, Le Roux N, Nienkemper-Swanepoel J, Ganey R, Van der Merwe C (2023) biplotEZ: EZ-to-Use Biplots. <https://CRAN.R-project.org/package=biplotEZ>, r package version 1.2.0
- McLachlan GJ (1992) Discriminant analysis and statistical pattern recognition. John Wiley, New York
- Rao CR (1948) The utilization of multiple measurements in problems of biological classification. *J Royal Stat Soc , Series B* 10:159–193
- da Rocha Neto AR, Sousa R, Barreto GdA, Cardoso JS (2011) Diagnostic of pathology on the vertebral column with embedded reject option. In: Iberian Conference on Pattern Recognition and Image Analysis, Springer, pp 588–595

- Ruts I, Rousseeuw PJ (1996) Computing depth contours of bivariate point clouds. *Comput Stat Data Anal* 23:153–168
- Tukey JW (1975) Mathematics and the picturing of data. *Proc Int Congress Math* 2:523–531

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.