

```

> X<-t(as.matrix(A[,1:p]) ## extract the data matrix X'
> one<-matrix(c(rep(1,n),n,1)
> nob<-c(n1,n2,...,nk) ## set up a vector of group sizes
> n<-sum(nob) ## total sample size
> G<-matrix(c(rep(1,nob[1]),rep(0,n),rep(1,nob[2])),
+ rep(0,n),...,rep(1,nk)),n,k) ## G is the group indicator
## matrix G_ij=1 if and only if case i is in group j
> F<-G%*%t(matrix(c(rep(1,nob[1])/nob[1],rep(0,n),
+ rep(1,nob[2])/nob[2],rep(0,n),...,rep(1,nk/nob[k])),n,k))
## F is used in calculating group means in M
> M<-F%*%t(X)
> xbar<-t(one)%*%t(X)/n ## overall mean vector
> W<-t(t(X)-G%*%M)%*%(t(X)-G%*%M)/(n-k) ##
> B<-t(G%*%M-one%*%xbar)%*%(G%*%M-one%*%xbar)/(k-1)
## W and B are within and between groups variances

```

A useful check on the calculations is to ensure that the analysis of variance is satisfied, i.e., that $(n-1) \cdot \text{var}(t(X)) = (n-k) \cdot W + (k-1) \cdot B$ (up to rounding errors). Note that the line calculating **B** may appear to be different from the formula given in the previous section $B = \frac{1}{k-1} \sum_{i=1}^k (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})'$ but the **R** calculations are in terms of matrices with n_i copies of $(\bar{x}_i - \bar{x})$.

If the number of groups k is large (i.e., more than 5 or 6, say) then calculation of the group indicator **G** is cumbersome and a quicker way of doing this **if the group sizes are equal** is to split the dataframe into separate group dataframes with the function `split(.)` (note that the argument of `split(.)` must be of class "dataframe" and not of class "matrix") followed by calculation of variances of each group dataframe with the function `lapply(.)` as follows:

```

Xdat<-split(A[,1:p],A[,p+1])
## note A is a dataframe not a matrix
## creates a list of dataframes, one for each group
Xdat<-lapply(Xdat,as.matrix)
## convert group dataframes to matrices
Xvar<-lapply(Xdat,var)
## find variance of each group
W<-Reduce("+",Xvar)*(n/k-1)/(n-k)
## assumes groups are of equal sizes
B<-((n-1)*var(X)-(W*(n-k))/k-1

```

Having obtained **W** and **B**, the within- and between-groups variances, it is then easy to calculate any of the statistics used for multivariate analysis of variance referred to in §9.2.5.2 and §9.2.7.3. The crimcoords can be obtained by the first $k-1$ vectors given by `eigen(solve(W)%*%B)$vectors` but note that the scaling used for the eigenvectors ensures that $y'y = 1$ whereas the scaling used for the crimcoords produced by the function `lda(.)` ensures that $y'Wy = 1$, i.e., multiplying the eigenvectors produced by `eigen(solve(W)%*%B)$vectors` by a

factor $(\mathbf{y}'\mathbf{y}/\mathbf{y}'\mathbf{W}\mathbf{y})^{1/2}$ will reproduce the crimcoords from `lda(.)` held in the matrix `lda(.)$scaling`.

The data transformed to crimcoords are given by $\mathbf{Y}' = (\mathbf{X} - \bar{\mathbf{X}})'\mathbf{V}$ where \mathbf{V} is obtained either from `V=eigen(solve(W)%*%B)$vectors[,1:k-1]` or from `V=lda(.)$scaling`. The difference in scaling will not be apparent if scatterplots of the data referred to crimcoords are produced by the basic **R** command `plot(.)` but will be noticeable if the MASS library routine `eqscplot(.)` is used.

9.5 Canonical Correlation Analysis

Canonical correlation analysis is concerned with investigating the relationship between two sets of variables measured on the same objects. In particular, the aim is to find which linear combination of variables of a $n \times p$ data matrix \mathbf{X}' has the maximum correlation with which linear combination of variables of a $n \times q$ data matrix \mathbf{Y}' amongst all such linear combinations. For example, in analyzing results of questionnaires eliciting subjects' opinions of a product, one set of variables may reflect the socio-economic aspects of the subjects and the other set may relate to their opinions on various properties of the product.

9.5.1 Derivation of canonical variates

If \mathbf{S}_{xx} , \mathbf{S}_{yy} and \mathbf{S}_{xy} are the sample variances and covariance matrices respectively of \mathbf{X}' , \mathbf{Y}' and $(\mathbf{X}', \mathbf{Y}')$ and \mathbf{x} and \mathbf{y} are p -vectors then the correlations between $\mathbf{X}'\mathbf{x}$ and $\mathbf{Y}'\mathbf{y}$ is $\mathbf{x}'\mathbf{S}_{xy}\mathbf{y}/\sqrt{\mathbf{x}'\mathbf{S}_{xx}\mathbf{x}\mathbf{y}'\mathbf{S}_{yy}\mathbf{y}}$. It was shown in Example 6 (iv) that this is maximized with respect to \mathbf{x} and \mathbf{y} by taking \mathbf{x} to be the eigenvector of $\mathbf{S}_{xx}^{-1}\mathbf{S}_{xy}\mathbf{S}_{yy}^{-1}\mathbf{S}'_{xy}$ corresponding to its largest eigenvalue and \mathbf{y} to be the eigenvector of $\mathbf{S}_{yy}^{-1}\mathbf{S}'_{xy}\mathbf{S}_{xx}^{-1}\mathbf{S}_{xy}$ corresponding to its largest eigenvalue. These eigenvectors are termed the first **canonical variates** of \mathbf{X}' and \mathbf{Y}' .

If the complete set of eigenpairs of $\mathbf{S}_{xx}^{-1}\mathbf{S}_{xy}\mathbf{S}_{yy}^{-1}\mathbf{S}'_{xy}$ is $(\mathbf{u}_1, \lambda_1), \dots, (\mathbf{u}_r, \lambda_r)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r$ and those of $\mathbf{S}_{yy}^{-1}\mathbf{S}'_{xy}\mathbf{S}_{xx}^{-1}\mathbf{S}_{xy}$ are $(\mathbf{v}_1, \lambda_1), \dots, (\mathbf{v}_r, \lambda_r)$ (noting the eigenvalues are identical) where $r = \min(p, q)$ then it can be shown that the linear combinations of \mathbf{X}' with \mathbf{Y}' given by $(\mathbf{u}_2, \mathbf{v}_2) \dots, (\mathbf{u}_r, \mathbf{v}_r)$ maximise the correlation between linear functions of the \mathbf{X}' and \mathbf{Y}' variables subject to the constraints of orthogonality with earlier ones and thus are termed the canonical variates of \mathbf{X}' and \mathbf{Y}' . Plots of the [mean corrected] data referred to canonical variates may provide insight into the structure of a relationship between the data sets, i.e., plots of $(\mathbf{X} - \bar{\mathbf{X}})'\mathbf{u}_i$ against $(\mathbf{X} - \bar{\mathbf{X}})'\mathbf{u}_j$ (typically with $j = i + 1$) or $(\mathbf{X} - \bar{\mathbf{X}})'\mathbf{u}_i$ against $(\mathbf{Y} - \bar{\mathbf{Y}})'\mathbf{v}_j$ and $(\mathbf{Y} - \bar{\mathbf{Y}})'\mathbf{v}_i$ against $(\mathbf{Y} - \bar{\mathbf{Y}})'\mathbf{v}_j$ can all be useful in informal investigation the structure of the data.

It can be shown that if the \mathbf{Y}' variables are group indicators or a set of binary dummy variables then the canonical variates of the \mathbf{X}' variables are precisely the