

Chapter 7

Canonical Correlation Analysis

7.1 Introduction

Let \mathbf{x} be the $p \times 1$ vector of predictors, and partition $\mathbf{x} = (\mathbf{w}^T, \mathbf{y}^T)^T$ where \mathbf{w} is $m \times 1$ and \mathbf{y} is $q \times 1$ where $m = p - q \leq q$ and $m, q \geq 2$. Canonical correlation analysis (CCA) seeks m pairs of linear combinations $(\mathbf{a}_1^T \mathbf{w}, \mathbf{b}_1^T \mathbf{y}), \dots, (\mathbf{a}_m^T \mathbf{w}, \mathbf{b}_m^T \mathbf{y})$ such that $\text{corr}(\mathbf{a}_i^T \mathbf{w}, \mathbf{b}_i^T \mathbf{y})$ is large under some constraints on the \mathbf{a}_i and \mathbf{b}_i where $i = 1, \dots, m$. The first pair $(\mathbf{a}_1^T \mathbf{w}, \mathbf{b}_1^T \mathbf{y})$ has the largest correlation. The next pair $(\mathbf{a}_2^T \mathbf{w}, \mathbf{b}_2^T \mathbf{y})$ has the largest correlation among all pairs uncorrelated with the first pair and the process continues so that $(\mathbf{a}_m^T \mathbf{w}, \mathbf{b}_m^T \mathbf{y})$ is the pair with the largest correlation that is uncorrelated with the first $m - 1$ pairs. The correlations are called *canonical correlations* while the pairs of linear combinations are called *canonical variables*.

Some notation is needed to explain CCA. Let the $p \times p$ positive definite symmetric dispersion matrix

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Let $\mathbf{J} = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}$. Let $\Sigma_a = \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$, $\Sigma_A = \mathbf{J} \mathbf{J}^T = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$, $\Sigma_b = \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ and $\Sigma_B = \mathbf{J}^T \mathbf{J} = \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$. Let \mathbf{e}_i and \mathbf{g}_i be sets of orthonormal eigenvectors, so $\mathbf{e}_i^T \mathbf{e}_i = 1$, $\mathbf{e}_i^T \mathbf{e}_j = 0$ for $i \neq j$, $\mathbf{g}_i^T \mathbf{g}_i = 1$ and $\mathbf{g}_i^T \mathbf{g}_j = 0$ for $i \neq j$. Let the \mathbf{e}_i be $m \times 1$ while the \mathbf{g}_i are $q \times 1$.

Let Σ_a have eigenvalue eigenvector pairs $(\lambda_1, \mathbf{a}_1), \dots, (\lambda_m, \mathbf{a}_m)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. Let Σ_A have eigenvalue eigenvector pairs $(\lambda_i, \mathbf{e}_i)$ for $i =$

$1, \dots, m$. Let Σ_b have eigenvalue eigenvector pairs $(\lambda_1, \mathbf{b}_1), \dots, (\lambda_q, \mathbf{b}_q)$. Let Σ_B have eigenvalue eigenvector pairs $(\lambda_i, \mathbf{g}_i)$ for $i = 1, \dots, q$. It can be shown that the m largest eigenvalues of the four matrices are the same. Hence $\lambda_i(\Sigma_a) = \lambda_i(\Sigma_A) = \lambda_i(\Sigma_b) = \lambda_i(\Sigma_B) \equiv \lambda_i$ for $i = 1, \dots, m$. It can be shown that $\mathbf{a}_i = \Sigma_{11}^{-1/2} \mathbf{e}_i$ and $\mathbf{b}_i = \Sigma_{22}^{-1/2} \mathbf{g}_i$. The eigenvectors \mathbf{a}_i are not necessarily orthonormal and the eigenvectors \mathbf{b}_i are not necessarily orthonormal.

Theorem 7.1. Assume the $p \times p$ dispersion matrix Σ is positive definite. Assume $\Sigma_{11}, \Sigma_{22}, \Sigma_A, \Sigma_a, \Sigma_B$ and Σ_b are positive definite and that $\hat{\Sigma} \xrightarrow{P} c\Sigma$ for some constant $c > 0$. Let \mathbf{d}_i be an eigenvector of the corresponding matrix. Hence $\mathbf{d}_i = \mathbf{a}_i, \mathbf{b}_i, \mathbf{e}_i$ or \mathbf{g}_i . Let $(\hat{\lambda}_i, \hat{\mathbf{d}}_i)$ be the i th eigenvalue eigenvector pair of $\hat{\Sigma}_\gamma$.

- a) $\hat{\Sigma}_\gamma \xrightarrow{P} \Sigma_\gamma$ and $\hat{\lambda}_i(\hat{\Sigma}_\gamma) \xrightarrow{P} \lambda_i(\Sigma_\gamma) = \lambda_i$ where $\gamma = A, a, B$ or b .
- b) $\Sigma_\gamma \hat{\mathbf{d}}_i - \lambda_i \hat{\mathbf{d}}_i \xrightarrow{P} \mathbf{0}$ and $\hat{\Sigma}_\gamma \mathbf{d}_i - \hat{\lambda}_i \mathbf{d}_i \xrightarrow{P} \mathbf{0}$.

c) If the j th eigenvalue λ_j is unique where $j \leq m$, then the absolute value of the correlation of $\hat{\mathbf{d}}_j$ with \mathbf{d}_j converges to 1 in probability: $|\text{corr}(\hat{\mathbf{d}}_j, \mathbf{d}_j)| \xrightarrow{P} 1$.

Proof. a) $\hat{\Sigma}_\gamma \xrightarrow{P} \Sigma_\gamma$ since matrix multiplication is a continuous function of the relevant matrices and matrix inversion is a continuous function of a positive definite matrix. Then $\hat{\lambda}_i(\hat{\Sigma}_\gamma) \xrightarrow{P} \lambda_i$ since an eigenvalue is a continuous function of its associated matrix.

b) Note that $(\Sigma_\gamma - \lambda_i \mathbf{I}) \hat{\mathbf{d}}_i = [(\Sigma_\gamma - \lambda_i \mathbf{I}) - (\hat{\Sigma}_\gamma - \hat{\lambda}_i \mathbf{I})] \hat{\mathbf{d}}_i = o_P(1) O_P(1) \xrightarrow{P} \mathbf{0}$, and $\hat{\Sigma}_\gamma \mathbf{d}_i - \hat{\lambda}_i \mathbf{d}_i \xrightarrow{P} \Sigma_\gamma \mathbf{d}_i - \lambda_i \mathbf{d}_i = \mathbf{0}$.

c) If n is large, then $\hat{\mathbf{d}}_i \equiv \hat{\mathbf{d}}_{i,n}$ is arbitrarily close to either \mathbf{d}_i or $-\mathbf{d}_i$, and the result follows.

Rule of thumb 7.1. To use CCA, assume the DD plot and subplots of the scatterplot matrix are linear. Want $n > 10p$ for classical CCA and $n > 20p$ for robust CCA that uses FCH, RFCH or RMVN. Also make the DD plot for the \mathbf{y} variables and the DD plot for the \mathbf{z} variables.

Definition 7.1. Let the dispersion matrix be $\text{Cov}(\mathbf{x}) = \Sigma_{\mathbf{x}}$. Let $(\lambda_i, \mathbf{e}_i)$ and $(\lambda_i, \mathbf{g}_i)$ be the eigenvalue eigenvector pairs of Σ_A and Σ_B . The k th pair of *population canonical variables* is

$$U_k = \mathbf{a}_k^T \mathbf{w} = \mathbf{e}_k^T \Sigma_{11}^{-1/2} \mathbf{w} \quad \text{and} \quad V_k = \mathbf{b}_k^T \mathbf{y} = \mathbf{g}_k^T \Sigma_{22}^{-1/2} \mathbf{y}$$

for $k = 1, \dots, m$. Then the *population canonical correlations* $\rho_k = \text{corr}(U_k, V_k)$

$= \sqrt{\lambda_k}$ for $k = 1, \dots, m$. The vectors $\mathbf{a}_k = \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_k$ and $\mathbf{b}_k = \boldsymbol{\Sigma}_{22}^{-1/2} \mathbf{g}_k$ are the k th *canonical correlation coefficient vectors* for \mathbf{w} and \mathbf{y} .

Theorem 7.2. Johnson and Wichern (1988, p. 440-441): Let the dispersion matrix be $\text{Cov}(\mathbf{x}) = \boldsymbol{\Sigma}_{\mathbf{x}}$. Then $V(U_k) = V(V_k) = 1$, $\text{Cov}(C_k, D_j) = \text{corr}(C_k, D_j) = 0$ for $k \neq j$ where $C_k = U_k$ or $C_k = V_k$, and $D_j = U_j$ or $D_j = V_j$ and $j, k = 1, \dots, m$. That is, U_k is uncorrelated with V_j and U_j for $j \neq k$, and V_k is uncorrelated with V_j and U_j for $j \neq k$. The first pair of canonical variables is the pair of linear combinations (U, V) having unit variances that maximizes $\text{corr}(U, V)$ and this maximum is $\text{corr}(U_1, V_1) = \rho_1$. The i th pair of canonical variables are the linear combinations (U, V) with unit variances that maximize $\text{corr}(U, V)$ among all choices uncorrelated with the previous $k - 1$ canonical variable pairs.

Definition 7.2. Suppose standardized data $\mathbf{z} = (\mathbf{w}^T, \mathbf{y}^T)^T$ is used and the dispersion matrix is the correlation matrix $\boldsymbol{\Sigma} = \boldsymbol{\rho}$. Hence $\boldsymbol{\Sigma}_{ii} = \boldsymbol{\rho}_{ii}$ for $i = 1, 2$. Let $(\lambda_i, \mathbf{e}_i)$ and $(\lambda_i, \mathbf{g}_i)$ be the eigenvalue eigenvector pairs of $\boldsymbol{\Sigma}_A$ and $\boldsymbol{\Sigma}_B$. The k th pair of *population canonical variables* is

$$U_k = \mathbf{a}_k^T \mathbf{w} = \mathbf{e}_k^T \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{w} \quad \text{and} \quad V_k = \mathbf{b}_k^T \mathbf{y} = \mathbf{g}_k^T \boldsymbol{\Sigma}_{22}^{-1/2} \mathbf{y}$$

for $k = 1, \dots, m$ for $k = 1, \dots, m$. Then the *population canonical correlations* $\rho_k = \text{corr}(U_k, V_k) = \sqrt{\lambda_k}$ for $k = 1, \dots, m$.

Then Theorem 7.2 holds for the standardized data and the canonical correlations are unchanged by the standardization.

Let

$$\hat{\boldsymbol{\Sigma}} = \begin{pmatrix} \hat{\boldsymbol{\Sigma}}_{11} & \hat{\boldsymbol{\Sigma}}_{12} \\ \hat{\boldsymbol{\Sigma}}_{21} & \hat{\boldsymbol{\Sigma}}_{22} \end{pmatrix}.$$

Define estimators $\hat{\boldsymbol{\Sigma}}_a$, $\hat{\boldsymbol{\Sigma}}_A$, $\hat{\boldsymbol{\Sigma}}_b$ and $\hat{\boldsymbol{\Sigma}}_B$ in the same manner as their population analogs but using $\hat{\boldsymbol{\Sigma}}$ instead of $\boldsymbol{\Sigma}$. For example, $\hat{\boldsymbol{\Sigma}}_a = \hat{\boldsymbol{\Sigma}}_{11}^{-1} \hat{\boldsymbol{\Sigma}}_{12} \hat{\boldsymbol{\Sigma}}_{22}^{-1} \hat{\boldsymbol{\Sigma}}_{21}$.

Let $\hat{\boldsymbol{\Sigma}}_a$ have eigenvalue eigenvector pairs $(\hat{\lambda}_i, \hat{\mathbf{a}}_i)$, and let $\hat{\boldsymbol{\Sigma}}_A$ have eigenvalue eigenvector pairs $(\hat{\lambda}_i, \hat{\mathbf{e}}_i)$ for $i = 1, \dots, m$. Let $\hat{\boldsymbol{\Sigma}}_b$ have eigenvalue eigenvector pairs $(\hat{\lambda}_1, \hat{\mathbf{b}}_1)$, and let $\hat{\boldsymbol{\Sigma}}_B$ have eigenvalue eigenvector pairs $(\hat{\lambda}_i, \hat{\mathbf{g}}_i)$ for $i = 1, \dots, q$. For these four matrices $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_m$.

Definition 7.3. Let $\hat{\boldsymbol{\Sigma}} = \mathbf{S}$ if data $\mathbf{x} = (\mathbf{w}^T, \mathbf{y}^T)^T$ is used, and let $\hat{\boldsymbol{\Sigma}} = \mathbf{R}$ if standardized data $\mathbf{z} = (\mathbf{w}^T, \mathbf{y}^T)^T$ is used. The k th pair of *sample*

canonical variables is

$$\hat{U}_k = \hat{\mathbf{a}}_k^T \mathbf{w} = \hat{\mathbf{e}}_k^T \hat{\boldsymbol{\Sigma}}_{11}^{-1/2} \mathbf{w} \quad \text{and} \quad \hat{V}_k = \hat{\mathbf{b}}_k^T \mathbf{y} = \hat{\mathbf{g}}_k^T \hat{\boldsymbol{\Sigma}}_{22}^{-1/2} \mathbf{y}$$

for $k = 1, \dots, m$. Then the *sample canonical correlations* $\hat{\rho}_k = \text{corr}(\hat{U}_k, \hat{V}_k) = \sqrt{\hat{\lambda}_k}$ for $k = 1, \dots, m$. The vectors $\hat{\mathbf{a}}_k = \hat{\boldsymbol{\Sigma}}_{11}^{-1/2} \hat{\mathbf{e}}_k$ and $\hat{\mathbf{b}}_k = \hat{\boldsymbol{\Sigma}}_{22}^{-1/2} \hat{\mathbf{g}}_k$ are the k th *sample canonical correlation vectors* for \mathbf{w} and \mathbf{y} .

Theorem 7.3. Under the conditions of Definition 7.3, the first pair of canonical variables (\hat{U}_1, \hat{V}_1) is the pair of linear combinations (\hat{U}, \hat{V}) having unit sample variances that maximizes the sample correlation $\text{corr}(\hat{U}, \hat{V})$ and this maximum is $\text{corr}(\hat{U}_1, \hat{V}_1) = \hat{\rho}_1$. The i th pair of canonical variables are the linear combinations (\hat{U}, \hat{V}) with unit sample variances that maximize the sample $\text{corr}(\hat{U}, \hat{V})$ among all choices uncorrelated with the previous $k-1$ canonical variable pairs.

7.2 Robust CCA

The *R* function `cancor` does classical CCA and the `mpack` function `rcancor` does robust CCA by applying `cancor` on the RMVN set: the subset of the data used to compute RMVN.

Some theory is simple: the FCH, RFCH and RMVN methods of RCCA produce consistent estimators of the k th canonical correlation ρ_k on a large class of elliptically contoured distributions.

To see this, suppose $\text{Cov}(\mathbf{x}) = c_{\mathbf{x}} \boldsymbol{\Sigma}$ and $\mathbf{C} \equiv \mathbf{C}(\mathbf{X}) \xrightarrow{P} c \boldsymbol{\Sigma}$ where $c_{\mathbf{x}} > 0$ and $c > 0$ are some constants. Then $\mathbf{C}_{XX}^{-1} \mathbf{C}_{XY} \mathbf{C}_{YY}^{-1} \mathbf{C}_{YX} \xrightarrow{P} \boldsymbol{\Sigma}_A = \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \boldsymbol{\Sigma}_{YY}^{-1} \boldsymbol{\Sigma}_{YX}$, and $\mathbf{C}_{YY}^{-1} \mathbf{C}_{YX} \mathbf{C}_{XX}^{-1} \mathbf{C}_{XY} \xrightarrow{P} \boldsymbol{\Sigma}_B = \boldsymbol{\Sigma}_{YY}^{-1} \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY}$. Note that $\boldsymbol{\Sigma}_A$ and $\boldsymbol{\Sigma}_B$ only depend on $\boldsymbol{\Sigma}$ and do not depend on the constants c or $c_{\mathbf{x}}$.

(If \mathbf{C} is also the classical covariance matrix applied to some subset of the data, then the correlation matrix $\mathbf{G} \equiv \mathbf{R}_{\mathbf{C}}$ applied to the same subset satisfies $\mathbf{G}_{XX}^{-1} \mathbf{G}_{XY} \mathbf{G}_{YY}^{-1} \mathbf{G}_{YX} \xrightarrow{P} \mathbf{R}_A = \mathbf{R}_{XX}^{-1} \mathbf{R}_{XY} \mathbf{R}_{YY}^{-1} \mathbf{R}_{YX}$, and $\mathbf{G}_{YY}^{-1} \mathbf{G}_{YX} \mathbf{G}_{XX}^{-1} \mathbf{G}_{XY} \xrightarrow{P} \mathbf{R}_B = \mathbf{R}_{YY}^{-1} \mathbf{R}_{YX} \mathbf{R}_{XX}^{-1} \mathbf{R}_{XY}$.)

Since eigenvalues are continuous functions of the associated matrix, and the FCH, RFCH and RMVN estimators are consistent estimators of $c_1 \boldsymbol{\Sigma}$, $c_2 \boldsymbol{\Sigma}$ and $c_3 \boldsymbol{\Sigma}$ on a large class of elliptically contoured distributions, Theorem

7.1 holds, so these three RCCA methods and `rcancor` produce consistent estimators the k th canonical correlation ρ_k on that class of distributions.

Example 7.1. Example 2.2 describes the mussel data. Log transformation were taken on *muscle mass* M , *shell width* W and on the *shell mass* S . Then x contained the two log mass measurements while y contains L , H and $\log(W)$. The robust and classical CCAs were similar, but the canonical coefficients were difficult to interpret since $\log(W)$ has different units than L and H . Hence the log transformation were taken on all five variables and output is shown below.

The data set zm contains x and y , and the DD plot showed case 48 was separated from the bulk of the data, but near the identity line. The DD plot for x showed two cases, 8 and 48, were separated from the bulk of the data. Also the plotted points did not cluster tightly about the identity line. The DD plot for y looked fine. The classical CCA produces output `$cor`, `$xcoef` and `$ycoef`. These are the canonical correlations, the \mathbf{a}_i and the \mathbf{b}_i . The labels for the RCCA are `outcor`, `outxcoef` and `outycoef`.

Note that the first correlation was about 0.98 while the second correlation was small. The RCCA is the CCA on the RMVN data set, which is contained in a compact ellipsoidal region. The variability of the truncated data set is less than that of the entire data set, hence expect the robust \mathbf{a}_i and \mathbf{b}_i to be larger in magnitude, ignoring sign, than that of the classical \mathbf{a}_i and \mathbf{b}_i , since the variance of each canonical variate is equal to one, and RCCA uses the truncated data. Note that \mathbf{a}_1 was roughly proportional to $\log(S)$ while \mathbf{b}_1 gave slightly higher weight for $\log(H)$ then $\log(W)$ and then $\log(L)$. Note that the five variables have high pairwise correlations, so $\log(M)$ was not important given that $\log(S)$ was in x . The second pair $(\mathbf{a}_2, \mathbf{b}_2)$ might be ignored since the second canonical correlation was very low.

```
> cancor(x,y)
$cor
[1] 0.9818605 0.1555381

$xcoef
      [,1]      [,2]
S 0.12650486 0.4077765
M 0.01897332 -0.4872522
```

```

$ycoef
[,1]      [,2]      [,3]
L 0.1567463  0.7277888  2.1935890
W 0.1605139  0.8650480 -1.0676419
H 0.2143781 -2.0634587 -0.8303862

$xcenter
S          M
4.563856 2.850187

$ycenter
L          W          H
5.472944 3.697654 4.723295

> rcancor(x,y)
$out
$out$cor
[1] 0.98596703 0.06797587

$out$xcoef
[,1]      [,2]
S 0.14966183 0.6460117
M 0.03236328 -0.8543387

$out$ycoef
[,1]      [,2]      [,3]
L 0.1625452  0.4237524 -2.8492678
W 0.2369692  1.5379681  0.9356495
H 0.2530324 -2.6806462  1.7785931

$out$xcenter
S          M
4.651941 2.948571

$out$ycenter
L          W          H
5.496255 3.728292 4.745839

```

7.3 Summary

1) Let \mathbf{x} be the $p \times 1$ vector of predictors, and partition $\mathbf{x} = (\mathbf{w}^T, \mathbf{y}^T)^T$ where \mathbf{w} is $m \times 1$ and \mathbf{y} is $q \times 1$ where $m = p - q \leq q$ and $m, q \geq 2$. Canonical correlation analysis (CCA) seeks m pairs of linear combinations $(\mathbf{a}_1^T \mathbf{w}, \mathbf{b}_1^T \mathbf{y}), \dots, (\mathbf{a}_m^T \mathbf{w}, \mathbf{b}_m^T \mathbf{y})$ such that $\text{corr}(\mathbf{a}_i^T \mathbf{w}, \mathbf{b}_i^T \mathbf{y})$ is large under some constraints on the \mathbf{a}_i and \mathbf{b}_i where $i = 1, \dots, m$. The first pair $(\mathbf{a}_1^T \mathbf{w}, \mathbf{b}_1^T \mathbf{y})$ has the largest correlation. The next pair $(\mathbf{a}_2^T \mathbf{w}, \mathbf{b}_2^T \mathbf{y})$ has the largest correlation among all pairs uncorrelated with the first pair and the process continues so that $(\mathbf{a}_m^T \mathbf{w}, \mathbf{b}_m^T \mathbf{y})$ is the pair with the largest correlation that is uncorrelated with the first $m - 1$ pairs. The correlations are called *canonical correlations* while the pairs of linear combinations are called *canonical variables*.

	corr		
	$\hat{\rho}_1$	\dots	$\hat{\rho}_1$
wcoef			
\mathbf{w}	$\hat{\mathbf{a}}_1$	\dots	$\hat{\mathbf{a}}_m$
ycoef			
\mathbf{y}	$\hat{\mathbf{b}}_1$	\dots	$\hat{\mathbf{b}}_m$
			$\hat{\mathbf{b}}_q$

```
64) $out$cor
[1] 0.98596703 0.06797587  $out$ycoef
$out$xcoef
[,1] [,2] [,3]
[1] 0.1625452 0.4237524 -2.8492678
S 0.14966183 0.6460117 W 0.2369692 1.5379681 0.9356495
M 0.03236328 -0.8543387 H 0.2530324 -2.6806462 1.7785931
```

3) Some notation is needed to explain CCA. Let the $p \times p$ positive definite symmetric dispersion matrix

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Let $\mathbf{J} = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}$. Let $\Sigma_a = \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$, $\Sigma_A = \mathbf{J} \mathbf{J}^T = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2}$, $\Sigma_b = \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$ and $\Sigma_B = \mathbf{J}^T \mathbf{J} = \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$. Let \mathbf{e}_i and \mathbf{g}_i be sets of orthonormal eigenvectors, so $\mathbf{e}_i^T \mathbf{e}_i = 1$, $\mathbf{e}_i^T \mathbf{e}_j = 0$ for $i \neq j$, $\mathbf{g}_i^T \mathbf{g}_i = 1$ and $\mathbf{g}_i^T \mathbf{g}_j = 0$ for $i \neq j$. Let the \mathbf{e}_i be $m \times 1$ while the \mathbf{g}_i are $q \times 1$.

Let Σ_a have eigenvalue eigenvector pairs $(\lambda_1, \mathbf{a}_1), \dots, (\lambda_m, \mathbf{a}_m)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. Let Σ_A have eigenvalue eigenvector pairs $(\lambda_i, \mathbf{e}_i)$ for $i =$

$1, \dots, m$. Let Σ_b have eigenvalue eigenvector pairs $(\lambda_1, \mathbf{b}_1), \dots, (\lambda_q, \mathbf{b}_q)$. Let Σ_B have eigenvalue eigenvector pairs $(\lambda_i, \mathbf{g}_i)$ for $i = 1, \dots, q$. It can be shown that the m largest eigenvalues of the four matrices are the same. Hence $\lambda_i(\Sigma_a) = \lambda_i(\Sigma_A) = \lambda_i(\Sigma_b) = \lambda_i(\Sigma_B) \equiv \lambda_i$ for $i = 1, \dots, m$. It can be shown that $\mathbf{a}_i = \Sigma_{11}^{-1/2} \mathbf{e}_i$ and $\mathbf{b}_i = \Sigma_{22}^{-1/2} \mathbf{g}_i$. The eigenvectors \mathbf{a}_i are not necessarily orthonormal and the eigenvectors \mathbf{b}_i are not necessarily orthonormal.

Theorem 7.1. Assume the $p \times p$ dispersion matrix Σ is positive definite. Assume $\Sigma_{11}, \Sigma_{22}, \Sigma_A, \Sigma_a, \Sigma_B$ and Σ_b are positive definite and that $\hat{\Sigma} \xrightarrow{P} c\Sigma$ for some constant $c > 0$. Let \mathbf{d}_i be an eigenvector of the corresponding matrix. Hence $\mathbf{d}_i = \mathbf{a}_i, \mathbf{b}_i, \mathbf{e}_i$ or \mathbf{g}_i . Let $(\hat{\lambda}_i, \hat{\mathbf{d}}_i)$ be the i th eigenvalue eigenvector pair of $\hat{\Sigma}_\gamma$.

- a) $\hat{\Sigma}_\gamma \xrightarrow{P} \Sigma_\gamma$ and $\hat{\lambda}_i(\hat{\Sigma}_\gamma) \xrightarrow{P} \lambda_i(\Sigma_\gamma) = \lambda_i$ where $\gamma = A, a, B$ or b .
- b) $\Sigma_\gamma \hat{\mathbf{d}}_i - \lambda_i \hat{\mathbf{d}}_i \xrightarrow{P} \mathbf{0}$ and $\hat{\Sigma}_\gamma \mathbf{d}_i - \hat{\lambda}_i \mathbf{d}_i \xrightarrow{P} \mathbf{0}$.
- c) If the j th eigenvalue λ_j is unique where $j \leq m$, then the absolute value of the correlation of $\hat{\mathbf{d}}_j$ with \mathbf{d}_j converges to 1 in probability: $|\text{corr}(\hat{\mathbf{d}}_j, \mathbf{d}_j)| \xrightarrow{P} 1$.

7.4 Complements

Muirhead and Waternaux (1980) shows that if the population canonical correlations ρ_k are distinct and if the underlying population distribution has a finite fourth moments, then the limiting joint distribution of $\sqrt{n}(\hat{\rho}_k^2 - \rho_k^2)$ is multivariate normal where the $\hat{\rho}_k$ are the classical sample canonical correlations and $k = 1, \dots, p$. If the data are iid from an elliptically contoured distribution with kurtosis 3κ , then the limiting joint distribution of

$$\sqrt{n} \frac{\hat{\rho}_k^2 - \rho_k^2}{2\rho_k(1 - \rho_k^2)}$$

for $k = 1, \dots, p$ is $N_p(\mathbf{0}, (\kappa + 1)\mathbf{I}_p)$. Note that $\kappa = 0$ for multivariate normal data.

Alkenani and Yu (2012), Zhang (2011) and Zhang, Olive and Ye (2012) develop robust CCA based on FCH, RFCH and RMVN.

7.5 Problems

PROBLEMS WITH AN ASTERISK * ARE ESPECIALLY USE-

FUL.

7.1*. Examine the *R* output in Example 7.1. a) What is the first canonical correlation $\hat{\rho}_1$?

- b) What is $\hat{\mathbf{a}}_1$?
- c) What is $\hat{\mathbf{b}}_1$?

7.2. The *R* output below is for a canonical correlation analysis on Venables and Ripley (2003) CPU data. The variables were $\text{syct} = \log(\text{cycle time} + 1)$,

$\text{mmin} = \log(\text{minimum main memory} + 1)$,

$\text{chmin} = \log(\text{minimum number of channels} + 1)$,

$\text{chmax} = \log(\text{maximum number of channels} + 1)$,

$\text{perf} = \log(\text{published performance} + 1)$ and

$\text{estperf} = 20/\sqrt{(\text{estimated performance} + 1)}$. These six variables had a linear scatterplot matrix and DD plot and similar variances. Want to compare the two performance variables with the four remaining variables.

- a) What is the first canonical correlation $\hat{\rho}_1$?

- b) What is $\hat{\mathbf{a}}_1$?

- c) What is $\hat{\mathbf{b}}_1$?

- d) Interpret the second canonical variable $U_2 = \hat{\mathbf{a}}_2^T \mathbf{w}$.

```
> cancor(w,y)
$cor
[1] 0.8769433 0.2278554

$xcoef
```

```

[,1]      [,2]
perf      0.02536432 0.1558717
estperf -0.04121870 0.1431100

$ycoef
[,1]      [,2]      [,3]      [,4]
syct   -0.013613254 0.05700360 0.089757416 -0.011423664
mmin    0.037485282 -0.01874858 0.084442460 0.005859654
chmin   0.006932264 0.09843612 -0.021782624 0.090756713
chmax   0.019998948 0.01159728 0.007855559 -0.094198608

```

7.3. Edited SAS output for SAS Institute (1985, p. 146) Fitness Club Data is given below for CCA. Three physiological and three exercise variables measured on 20 middle aged men at a fitness club.

a) What is the first canonical correlation $\hat{\rho}_1$?

b) What is \hat{a}_1 ?

c) What is \hat{b}_1 ?

```

Canonical
Correlation
0.7956
0.2006
0.0726

```

Raw Canonical Coefficients for the Physiological Variables

	PHYS1	PHYS2	PHYS3
weight	-0.0314	-0.0763	-0.0077
waist	0.0493	0.3687	0.1580
pulse	-0.0082	-0.0321	0.1457

Raw Canonical Coefficients for the Exercise Variables

	Exer1	Exer2	Exer3
chinups	-0.0661	-0.0714	-0.2428
situps	-0.0168	0.0020	0.0198
jumps	0.0140	0.0207	-0.0082

7.4. The output below is for a canonical correlations analysis on the *R* Seatbelts data set where $y_1 = \text{drivers}$ = number of drivers killed or seriously injured, $y_2 = \text{front}$ = number of front seat passengers killed or seriously injured, and $y_3 = \text{rear}$ = number of back seat passengers killed or seriously injured, $x_1 = \text{kms}$ = distance driven, $x_2 = \text{PetrolPrice}$ = petrol price and $x_3 = \text{VanKilled}$ = number of van drivers killed. The data consists of 192 monthly totals in Great Britain from January 1969 to December 1984.

- a) What is the first canonical correlation $\hat{\rho}_1$?
- b) What is $\hat{\mathbf{a}}_1$?
- c) What is $\hat{\mathbf{b}}_1$?
- d) Let $\mathbf{z} = (\mathbf{x}^T, \mathbf{y}^T)^T$. From the DD plot, the \mathbf{z}_i appeared to follow a multivariate normal distribution. Sketch the DD plot.

```
> rcancor(x,y)
$out
$out$cor
```

```
[1] 0.8116953 0.5064619 0.1376399
```

```
$out$xcoef
```

	[,1]	[,2]	[,3]
x.kms	-2.080206e-05	-0.0000233873	-2.259723e-06
x.PetrolPrice	-1.847967e+00	3.7173715818	5.292041e+00
x.VanKilled	1.597620e-03	-0.0168450843	1.673662e-02

```
$out$ycoef
```

	[,1]	[,2]	[,3]
y.drivers	1.678751e-06	-2.487259e-05	0.0004717902
y.front	5.594715e-04	-7.797027e-05	-0.0008157585
y.rear	-9.964980e-04	-7.521578e-04	0.0005045756

7.5. The *R* output below is for a canonical correlation analysis on some iris data. An iris is a flower, and there were 50 observations with 4 variables sepal length, sepal width, petal length and petal width.

a) What is the first canonical correlation $\hat{\rho}_1$?

b) What is \hat{a}_1 ?

c) What is \hat{b}_1 ?

```
w<-iris3[, ,3]
x <- w[, 1:2]
y <- w[, 3:4]
cancor(x, y)
```

```
$cor
[1] 0.8642869 0.4836991
```

```
$xcoef
```

	[,1]	[,2]
Sepal L.	-0.223034210	-0.1186117

Sepal W. -0.006920448 0.4980378

\$ycoef

[,1] [,2]

Petal L. -0.257853414 -0.09094352

Petal W. -0.006108292 0.54939125