

8 Principal Components Analysis

Further reading. The paper Xoplaki et al. (2000) shows a nice application of the principal components and canonical correlation analysis method (see Chapter 9).

8.1 Introduction

Example (Weather Report, p. 1-6). The data set consists of sunshine (three variables), air temperature (six variables), heating degree days (HDD) and precipitation data (five variables). Figure 8.1 shows a scatterplot matrix of several variables. With the help of principal component analysis we want to reduce the number of variables, without losing a lot of information.

A principal component analysis is concerned with explaining the variance-covariance structure of a set of variables through a few linear combinations of these variables. Its general objectives are (1) data reduction and (2) interpretation.

Although p components are required to reproduce the total system variability, often much of this variability can be accounted for by a small number k of the principal components. If so, there is almost as much information in the k components as there is in the original p variables. The k principal components can then replace the initial p variables, and the original data set, consisting of n measurements on p variables, is reduced to a data set consisting of n measurements on k principal components.

Analyses of principal components are more of a means to an end rather than an end in themselves, because they frequently serve as intermediate steps in much larger investigations.

8.2 Population Principal Components

Algebraically, principal components are particular linear combinations of the p random variables X_1, \dots, X_p . Geometrically, these linear combinations represent the selection of a new coordinate system obtained by rotating the original system with X_1, \dots, X_p as the coordinate axes. The new axes represent the directions with maximum variability and provide a simpler and more parsimonious description of the covariance structure.

As we shall see, principal components depend solely on the covariance matrix Σ (or the correlation matrix ρ of X_1, \dots, X_p). Their development does not require a multivariate normal assumption. On the other hand, principal components derived for multivariate normal populations have useful interpretations in terms of the constant density ellipsoids. Further, inferences can be done from the sample components when the population is multivariate normal.

Let the random vector $\mathbf{X}' = (X_1, \dots, X_p)$ have the covariance matrix Σ with eigen-

values $\lambda_1 \geq \dots \geq \lambda_p \geq 0$. Consider the linear combinations

$$\begin{aligned} Y_1 &= \mathbf{a}'_1 \mathbf{X} = a_{11}X_1 + \dots + a_{1p}X_p \\ &\vdots \\ Y_p &= \mathbf{a}'_p \mathbf{X} = a_{p1}X_1 + \dots + a_{pp}X_p. \end{aligned}$$

Then, we find

$$\text{Var}(Y_i) = \mathbf{a}'_i \boldsymbol{\Sigma} \mathbf{a}_i, \quad i = 1, \dots, p \quad (8.1)$$

$$\text{Cov}(Y_i, Y_k) = \mathbf{a}'_i \boldsymbol{\Sigma} \mathbf{a}_k, \quad i, k = 1, \dots, p \quad (8.2)$$

The principal components are those uncorrelated linear combinations Y_1, \dots, Y_p whose variances in (8.1) are as large as possible.

Definition 8.2.1. We define

- First principal component = linear combination $\mathbf{a}'_1 \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_1 \mathbf{X})$ subject to $\mathbf{a}'_1 \mathbf{a}_1 = 1$
- Second principal component = linear combination $\mathbf{a}'_2 \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_2 \mathbf{X})$ subject to $\mathbf{a}'_2 \mathbf{a}_2 = 1$ and $\text{Cov}(\mathbf{a}'_1 \mathbf{X}, \mathbf{a}'_2 \mathbf{X}) = 0$
- \vdots
- i th principal component = linear combination $\mathbf{a}'_i \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_i \mathbf{X})$ subject to $\mathbf{a}'_i \mathbf{a}_i = 1$ and $\text{Cov}(\mathbf{a}'_i \mathbf{X}, \mathbf{a}'_k \mathbf{X}) = 0$ for $k < i$.

Proposition 8.2.2. Let $\boldsymbol{\Sigma}$ be the covariance matrix associated with the random vector $\mathbf{X}' = (X_1, \dots, X_p)$. Let $\boldsymbol{\Sigma}$ have eigenvalue-eigenvector pairs $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$ where $\lambda_1 \geq \dots \geq \lambda_p \geq 0$. Then the i th principal component is given by

$$Y_i = \mathbf{e}'_i \mathbf{X} = \sum_{j=1}^p e_{ij} X_j, \quad i = 1, \dots, p. \quad (8.3)$$

With these choices

$$\begin{aligned} \text{Var}(Y_i) &= \mathbf{e}'_i \boldsymbol{\Sigma} \mathbf{e}_i = \lambda_i, & i = 1, \dots, p \\ \text{Cov}(Y_i, Y_k) &= \mathbf{e}'_i \boldsymbol{\Sigma} \mathbf{e}_k = 0, & i \neq k. \end{aligned}$$

Remark. If some λ_i are equal, the choices of the corresponding coefficient vectors, \mathbf{e}_i , and hence Y_i , are not unique.

Proposition 8.2.3. Let $\mathbf{X}' = (X_1, \dots, X_p)$ as in Proposition 8.2.2. Then

$$\text{tr}(\boldsymbol{\Sigma}) = \sum_{i=1}^p \sigma_{ii} = \sum_{i=1}^p \text{Var}(X_i) = \sum_{i=1}^p \lambda_i = \sum_{i=1}^p \text{Var}(Y_i) \quad (8.4)$$

and the proportion of total population variance due to the k th principal component

$$\frac{\lambda_k}{\sum_{i=1}^p \lambda_i}, \quad k = 1, \dots, p.$$

Remark. If most of the total population variance, for large p , can be attributed to the first one, two or three components, then these components can “replace” the original p variables without much loss of information.

Remark. The magnitude of e_{ik} – also principal component loading – measures the importance of the k th variable to the i th principal component and thus is a useful basis for interpretation. A large coefficient (in absolute value) corresponds to a high loading, while a coefficient near zero has a low loading.

Remark. One important use of principal components is interpreting the original data in terms of the principal components. The images of the original data under the principal components transformation are referred to as principal component scores.

Proposition 8.2.4. If $Y_1 = \mathbf{e}'_1 \mathbf{X}, \dots, Y_p = \mathbf{e}'_p \mathbf{X}$ are the principal components obtained from the covariance matrix $\boldsymbol{\Sigma}$ then

$$\rho_{Y_i, X_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, \quad i, k = 1, \dots, p.$$

Remark. Suppose $\mathbf{X} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma})$. Then we have $c^2 = \sum_{i=1}^p y_i^2 / \lambda_i$ and this equation defines an ellipsoid with axes y_1, \dots, y_p lying in directions of $\mathbf{e}_1, \dots, \mathbf{e}_p$, respectively.

8.2.1 Principal Components obtained from Standardized Values

Principal components may also be obtained for the standardized variables

$$Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}, \quad i = 1, \dots, p \quad (8.5)$$

or

$$\mathbf{Z} = (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu}), \quad (8.6)$$

where the diagonal standard deviation matrix $\mathbf{V}^{1/2}$ is defined as

$$\mathbf{V}^{1/2} = \begin{pmatrix} \sqrt{\sigma_{11}} & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_{22}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_{pp}} \end{pmatrix}.$$

We find $E(\mathbf{Z}) = \mathbf{0}$ and $\text{Cov}(\mathbf{Z}) = (\mathbf{V}^{-1/2})^{-1} \Sigma (\mathbf{V}^{1/2})^{-1} = \boldsymbol{\rho}$.

Remark. The eigenvalue-eigenvector pairs derived from Σ are in general not the same as the ones derived from $\boldsymbol{\rho}$.

Proposition 8.2.5. *The i th principal component of the standardized variables $\mathbf{Z}' = (Z_1, \dots, Z_p)$ with $\text{Cov}(\mathbf{Z}) = \boldsymbol{\rho}$, is given by*

$$Y_i = \mathbf{e}_i' \mathbf{Z} = \mathbf{e}_i' (\mathbf{V}^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu}), \quad i = 1, \dots, p.$$

Moreover

$$\sum_{i=1}^p \text{Var}(Y_i) = \sum_{i=1}^p \text{Var}(Z_i) = p$$

and

$$\rho_{Y_i, Z_k} = e_{ik} \sqrt{\lambda_i}, \quad i, k = 1, \dots, p.$$

In this case $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$ are the eigenvalue-eigenvector pairs of $\boldsymbol{\rho}$ with $\lambda_1 \geq \dots \geq \lambda_p \geq 0$.

Remark. Variables should be standardized if they are measured on scales with widely differing ranges or if the units of measurements are not commensurate.

8.3 Summarizing Sample Variation by Principal Components

We now study the problem of summarizing the variation in n measurements on p variables with a few judiciously chosen linear combinations. Suppose the data $\mathbf{x}_1, \dots, \mathbf{x}_n$ represent n independent drawings from some p -dimensional population with mean vector $\boldsymbol{\mu}$ and covariance matrix Σ . These data yield the sample mean vector $\bar{\mathbf{x}}$, sample covariance matrix \mathbf{S} and the sample correlation matrix \mathbf{R} .

The uncorrelated combinations with the largest variances will be called the sample principal components. The sample principal components (PC) are defined as those linear

combinations which have maximum sample variance. Specifically,

First sample PC = linear combination $\mathbf{a}'_1 \mathbf{x}_j$ that maximizes the sample variance $\text{Var}(\mathbf{a}'_1 \mathbf{x}_j)$ subject to $\mathbf{a}'_1 \mathbf{a}_1 = 1$

Second sample PC = linear combination $\mathbf{a}'_2 \mathbf{x}_j$ that maximizes the sample variance $\text{Var}(\mathbf{a}'_2 \mathbf{x}_j)$ subject to $\mathbf{a}'_2 \mathbf{a}_2 = 1$ and $\text{Cov}(\mathbf{a}'_1 \mathbf{x}_j, \mathbf{a}'_2 \mathbf{x}_j) = 0$

\vdots

i th sample PC = linear combination $\mathbf{a}'_i \mathbf{x}_j$ that maximizes the sample variance $\text{Var}(\mathbf{a}'_i \mathbf{x}_j)$ subject to $\mathbf{a}'_i \mathbf{a}_i = 1$ and $\text{Cov}(\mathbf{a}'_i \mathbf{x}_j, \mathbf{a}'_k \mathbf{x}_j) = 0, \quad k < i.$

Proposition 8.3.1. *If $\mathbf{S} = \{s_{ik}\}$ is the $p \times p$ sample covariance matrix with eigenvalue-eigenvector pairs $(\hat{\lambda}_1, \hat{\mathbf{e}}_1), \dots, (\hat{\lambda}_p, \hat{\mathbf{e}}_p)$ the i th sample principal component is given by*

$$\hat{y}_i = \hat{\mathbf{e}}'_i \mathbf{x} = \hat{e}_{i1}x_1 + \dots + \hat{e}_{ip}x_p, \quad i = 1, \dots, p,$$

where $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p \geq 0$ and \mathbf{x} is any observation on the variable X_1, \dots, X_p . Also,

$$\begin{aligned} \text{sample variance } (\hat{y}_k) &= \hat{\lambda}_k, & k = 1, \dots, p, \\ \text{sample covariance } (\hat{y}_i, \hat{y}_k) &= 0, & i \neq k. \end{aligned}$$

In addition

$$\text{total sample variance} = \sum_{i=1}^p s_{ii} = \hat{\lambda}_1 + \dots + \hat{\lambda}_p$$

and

$$r(\hat{y}_i, x_k) = \frac{\hat{e}_{ik} \sqrt{\hat{\lambda}_i}}{\sqrt{\sigma_{kk}}}, \quad i, k = 1, \dots, p.$$

Number of Principal Components

How many components to retain? There is no definitive answer to this question. A useful visual aid to determining an appropriate number of principal components is a scree plot. With the eigenvalues ordered from largest to smallest, a scree plot is a plot of $\hat{\lambda}_i$ versus i . To determine the appropriate number of components, we look for an elbow (bend) in the scree plot. The number of components is taken to be the point at which the remaining eigenvalues are relatively small and all about the same size.

Remark. An unusually small value for the last eigenvalue from either the sample covariance or correlation matrix can indicate an unnoticed linear dependency in the data set and should therefore not be routinely ignored.

Example (Weather Report, p. 1-6). Figure 8.2 shows the results of the principal component analysis calculated with the correlation matrix.

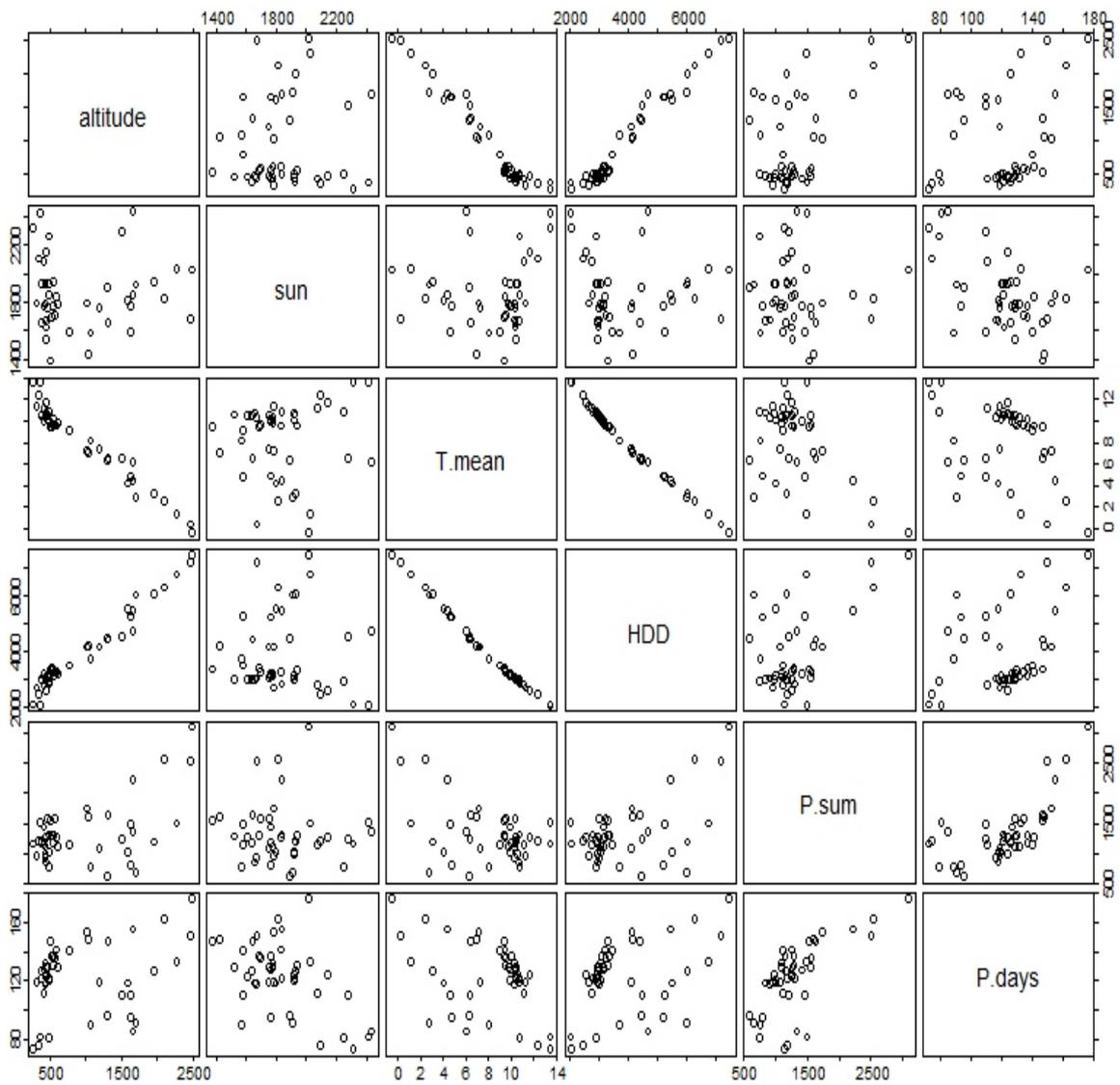


Figure 8.1: Scatterplot matrix of some variables of the Weather Report, p. 1-6.

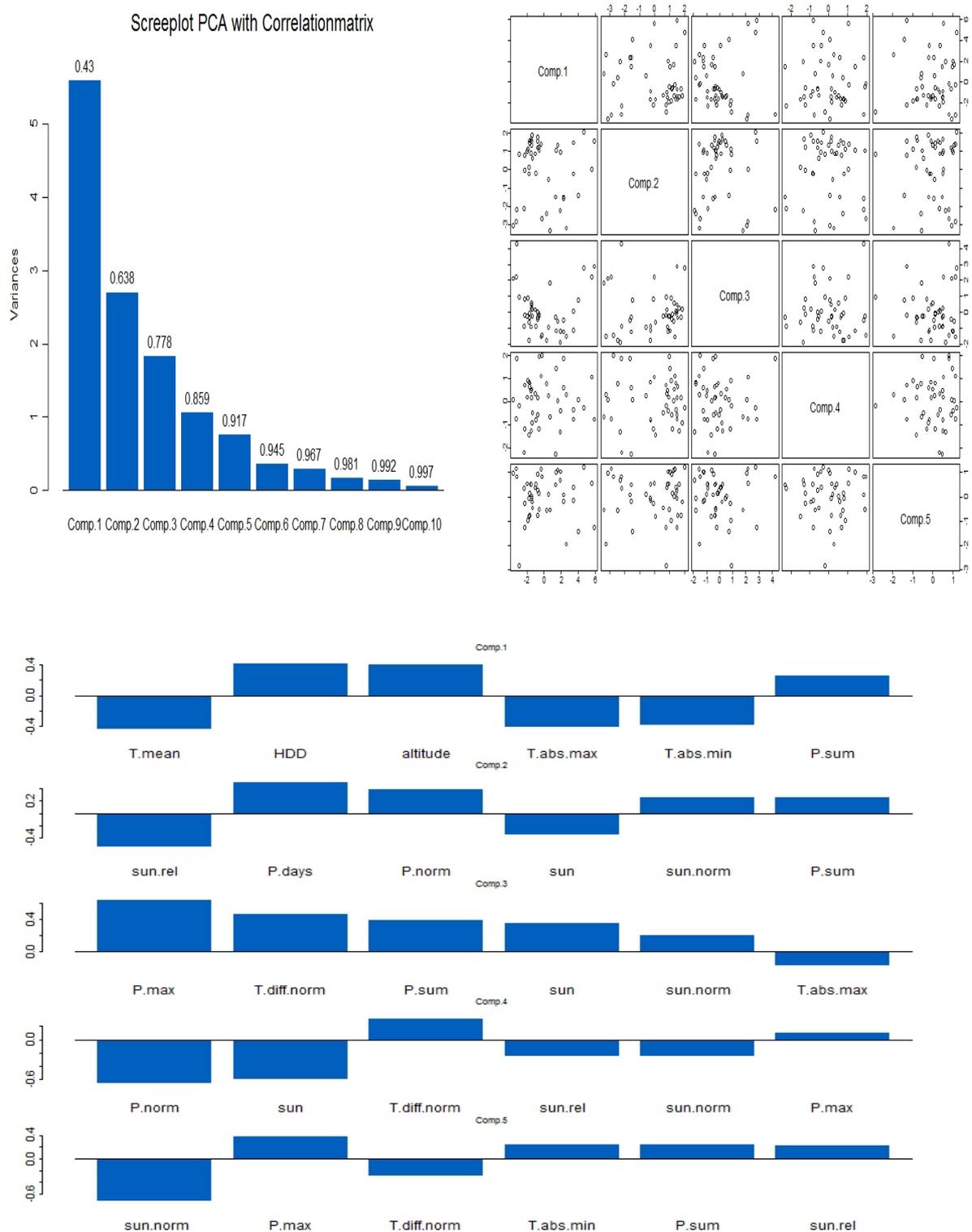


Figure 8.2: Screeplot of the principal components analysis (top left), scatterplot matrix of the scores calculated with the correlation matrix (top right) and the loadings of the variables (bottom). Data set: Weather Report, p. 1-6.

9 Canonical Correlation Analysis

Further reading. Read again the paper Xoplaki et al. (2000) to see how the canonical correlation analysis method can be applied in climate sciences.

9.1 Introduction

Example (Soil Evaporation, p. 1-8). The observed variables are maximum (maxst), minimum (minst), and average soil temperature (avst); maximum (maxat), minimum (minat), and average air temperature (avat); maximum (maxh), minimum (minh), and average relative humidity (avh); total wind in miles per day (wind) and the daily amount of evaporation from the soil (evap). The three “average” measurements are integrated: average soil temperature is the integrated area under the daily soil temperature curve, average air temperature is the integrated area under the daily air temperature curve, and average relative humidity is the integrated area under the daily relative humidity curve.

We want to find the association between the soil variables (maxst, minst and avst) and the air variables (maxat, minat, avat, maxh, minh, avh, wind).

Canonical correlation analysis (CCA) seeks to identify and quantify the associations between two sets of variables. Canonical correlation analysis focuses on the correlation between a linear combination of the variables in one set and a linear combination of the variables in another set. The idea is first to determine the pair of linear combinations having the largest correlation. Next, we determine the pair of linear combinations having the largest correlation among all pairs uncorrelated with the initially selected pair, and so on. The pairs of linear combinations are called the canonical variables, and their correlations are called canonical correlations. The canonical correlations measure the strength of association between the two sets of variables. The maximization aspect of the technique represents an attempt to concentrate a high-dimensional relationship between two sets of variables into a few pairs of canonical variables.

9.2 Canonical Variates and Canonical Correlations

We are interested in measures of association between two groups of variables. The first group, of p variables, is represented by the $(p \times 1)$ random vector $\mathbf{X}^{(1)}$. The second group, of q variables, is represented by the $(q \times 1)$ random vector $\mathbf{X}^{(2)}$. We assume, in the theoretical development, that $\mathbf{X}^{(1)}$ represents the smaller set, so that $p \leq q$.

For the random vector $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$, let

$$\begin{aligned} E(\mathbf{X}^{(1)}) &= \boldsymbol{\mu}^{(1)}, & \text{Cov}(\mathbf{X}^{(1)}) &= \boldsymbol{\Sigma}_{11} \\ E(\mathbf{X}^{(2)}) &= \boldsymbol{\mu}^{(2)}, & \text{Cov}(\mathbf{X}^{(2)}) &= \boldsymbol{\Sigma}_{22} \\ \text{Cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) &= \boldsymbol{\Sigma}_{12} = \boldsymbol{\Sigma}'_{21}. \end{aligned}$$

It will be convenient to consider $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ jointly, so we find that the random vector

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}^{(1)} \\ \mathbf{X}^{(2)} \end{pmatrix}$$

$((p+q) \times 1)$

has mean vector

$$\boldsymbol{\mu} = E(\mathbf{X}) = \begin{pmatrix} \boldsymbol{\mu}^{(1)} \\ \boldsymbol{\mu}^{(2)} \end{pmatrix}$$

$((p+q) \times 1)$

and covariance matrix

$$\begin{aligned} \boldsymbol{\Sigma} &= E(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})' \\ &= \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \begin{matrix} (p \times (p+q)) \\ (q \times (p+q)) \end{matrix} \\ &\quad \begin{matrix} ((p+q) \times p) & ((p+q) \times q) \end{matrix} \end{aligned}$$

The covariances between pairs of variables from different sets – one variable from $\mathbf{X}^{(1)}$, one variable from $\mathbf{X}^{(2)}$ – are contained in $\boldsymbol{\Sigma}_{12}$ or, equivalently, in $\boldsymbol{\Sigma}_{21}$. That is, the pq elements of $\boldsymbol{\Sigma}_{12}$ measure the association between the two sets. When p and q are relatively large, interpreting the elements of $\boldsymbol{\Sigma}_{12}$ collectively is ordinarily hopeless. Moreover, it is often linear combinations of variables that are interesting and useful for predictive or comparative purposes. The main task of canonical correlation analysis is to summarize the associations between the $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ sets in terms of a few carefully chosen covariances (or correlations) rather than the pq covariances in $\boldsymbol{\Sigma}_{12}$.

Linear combinations provide simple summary measures of a set of variables. Set

$$\begin{aligned} U &= \mathbf{a}'\mathbf{X}^{(1)} \\ V &= \mathbf{b}'\mathbf{X}^{(2)} \end{aligned}$$

for some pair of coefficient vectors \mathbf{a} and \mathbf{b} . We obtain

$$\begin{aligned} \text{Var}(U) &= \mathbf{a}'\text{Cov}(\mathbf{X}^{(1)})\mathbf{a} = \mathbf{a}'\boldsymbol{\Sigma}_{11}\mathbf{a} \\ \text{Var}(V) &= \mathbf{b}'\text{Cov}(\mathbf{X}^{(2)})\mathbf{b} = \mathbf{b}'\boldsymbol{\Sigma}_{22}\mathbf{b} \\ \text{Cov}(U, V) &= \mathbf{a}'\text{Cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)})\mathbf{b} = \mathbf{a}'\boldsymbol{\Sigma}_{12}\mathbf{b}. \end{aligned}$$

We shall seek coefficient vectors \mathbf{a} and \mathbf{b} such that

$$\text{Cor}(U, V) = \frac{\mathbf{a}'\boldsymbol{\Sigma}_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\boldsymbol{\Sigma}_{11}\mathbf{a}}\sqrt{\mathbf{b}'\boldsymbol{\Sigma}_{22}\mathbf{b}}} \quad (9.1)$$

is as large as possible.

Definition 9.2.1. We define:

First pair of canonical variables (first canonical variate pair): pair of linear combinations U_1, V_1 having unit variances, which maximizes the correlation (9.1)

Second pair of canonical variables (second canonical variate pair): pair of linear combinations U_2, V_2 having unit variances, which maximize the correlation (9.1) among all choices that are uncorrelated with the first pair of canonical variables.

At the k th step,

The k th pair of canonical variables (k th canonical variate pair): pair of linear combinations U_k, V_k having unit variances, which maximize the correlation (9.1) among all choices uncorrelated with the previous $k - 1$ canonical variable pairs.

The correlation between the k th pair of canonical variables is called the k th canonical correlation.

Proposition 9.2.2. Suppose $p \leq q$ and let the random vectors $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ have $\text{Cov}(\mathbf{X}^{(1)}) = \boldsymbol{\Sigma}_{11}$, $\text{Cov}(\mathbf{X}^{(2)}) = \boldsymbol{\Sigma}_{22}$ and $\text{Cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) = \boldsymbol{\Sigma}_{12} = \boldsymbol{\Sigma}'_{21}$, where $\boldsymbol{\Sigma}$ has full rank. For coefficient vectors \mathbf{a} and \mathbf{b} , form the linear combinations $U = \mathbf{a}'\mathbf{X}^{(1)}$ and $V = \mathbf{b}'\mathbf{X}^{(2)}$. Then

$$\max_{\mathbf{a}, \mathbf{b}} \text{Cor}(U, V) = \rho_1^*$$

is attained by the linear combinations

$$U_1 = \underbrace{\mathbf{e}'_1 \boldsymbol{\Sigma}_{11}^{-1/2}}_{\mathbf{a}'_1} \mathbf{X}^{(1)} \quad \text{and} \quad V_1 = \underbrace{\mathbf{f}'_1 \boldsymbol{\Sigma}_{22}^{-1/2}}_{\mathbf{b}'_1} \mathbf{X}^{(2)}.$$

The k th pair of canonical variates, $k = 2, \dots, p$,

$$U_k = \mathbf{e}'_k \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{X}^{(1)} \quad \text{and} \quad V_k = \mathbf{f}'_k \boldsymbol{\Sigma}_{22}^{-1/2} \mathbf{X}^{(2)}$$

maximizes

$$\text{Cor}(U_k, V_k) = \rho_k^*$$

among those linear combinations uncorrelated with the preceding $1, 2, \dots, k - 1$ canonical variables.

Here $(\rho_1^*)^2 \geq \dots \geq (\rho_p^*)^2$ are the eigenvalues of $\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2}$, and $\mathbf{e}_1, \dots, \mathbf{e}_p$ are the associated $(p \times 1)$ eigenvectors.

The canonical variates have the properties

$$\begin{aligned} \text{Var}(U_k) &= \text{Var}(V_k) = 1 \\ \text{Cov}(U_k, U_l) &= \text{Cov}(V_k, V_l) = \text{Cov}(U_k, V_l) = 0 \\ \text{Cor}(U_k, U_l) &= \text{Cor}(V_k, V_l) = \text{Cor}(U_k, V_l) = 0. \end{aligned}$$

for $k, l = 1, 2, \dots, p$ with $k \neq l$.

Remark. The quantities $(\rho_1^*)^2 \geq \dots \geq (\rho_p^*)^2$ are also the p largest eigenvalues of the matrix

$$\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1/2}$$

with corresponding $(q \times 1)$ eigenvectors $\mathbf{f}_1, \dots, \mathbf{f}_p$. Each \mathbf{f}_i is proportional to

$$\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_i.$$

Remark. If the original variables are standardized with $\mathbf{Z}^{(1)} = (Z_1^{(1)}, \dots, Z_p^{(1)})'$ and $\mathbf{Z}^{(2)} = (Z_1^{(2)}, \dots, Z_q^{(2)})'$, the canonical variates are of the form

$$\begin{aligned} U_k &= \mathbf{a}'_k \mathbf{Z}^{(1)} = \mathbf{e}'_k \boldsymbol{\rho}_{11}^{-1/2} \mathbf{Z}^{(1)} \\ V_k &= \mathbf{b}'_k \mathbf{Z}^{(2)} = \mathbf{f}'_k \boldsymbol{\rho}_{22}^{-1/2} \mathbf{Z}^{(2)}. \end{aligned} \quad (9.2)$$

Here, $\text{Cov}(\mathbf{Z}^{(1)}) = \boldsymbol{\rho}_{11}$, $\text{Cov}(\mathbf{Z}^{(2)}) = \boldsymbol{\rho}_{22}$, $\text{Cov}(\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}) = \boldsymbol{\rho}_{12} = \boldsymbol{\rho}'_{21}$, and \mathbf{e}_k and \mathbf{f}_k are the eigenvectors of $\boldsymbol{\rho}_{11}^{-1/2} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1} \boldsymbol{\rho}_{21} \boldsymbol{\rho}_{11}^{-1/2}$ and $\boldsymbol{\rho}_{22}^{-1/2} \boldsymbol{\rho}_{21} \boldsymbol{\rho}_{11}^{-1} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1/2}$, respectively. The canonical correlations, ρ_k^* , satisfy

$$\text{Cor}(U_k, V_k) = \rho_k^*, \quad k = 1, \dots, p,$$

where $(\rho_1^*)^2 \geq \dots \geq (\rho_p^*)^2$ are the nonzero eigenvalues of the matrix $\boldsymbol{\rho}_{11}^{-1/2} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1} \boldsymbol{\rho}_{21} \boldsymbol{\rho}_{11}^{-1/2}$ or, equivalently, the largest eigenvalues of $\boldsymbol{\rho}_{22}^{-1/2} \boldsymbol{\rho}_{21} \boldsymbol{\rho}_{11}^{-1} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1/2}$.

Remark. The canonical coefficients for the standardized variables,

$$Z_i^{(1)} = \frac{X_i^{(1)} - \mu_i^{(1)}}{\sqrt{\sigma_{ii}}},$$

are simply related to the canonical coefficients attached to the original variables $X_i^{(1)}$. Specifically, if \mathbf{a}'_k is the coefficient vector for the k th canonical variate U_k , then $\mathbf{a}'_k \mathbf{V}_{11}^{1/2}$ is the coefficient vector for the k th canonical variate constructed from the standardized variables $\mathbf{Z}^{(1)}$. Here $\mathbf{V}_{11}^{1/2}$ is the diagonal matrix with the i th diagonal element $\sqrt{\sigma_{ii}} = \sqrt{\text{Var}(X_i^{(1)})}$. Similarly, $\mathbf{b}'_k \mathbf{V}_{22}^{1/2}$ is the coefficient vector for the canonical variate constructed from the set of standardized variables $\mathbf{Z}^{(2)}$. The canonical correlations are unchanged by the standardization.

9.3 Interpreting the Population Canonical Variables

Canonical variables are, in general, artificial. That is, they have no physical meaning. If the original variables $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ are used, the canonical coefficients \mathbf{a} and \mathbf{b} have units proportional to those of the $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ sets. If the original variables are standardized to have zero means and unit variances, the canonical coefficients have no units of measurement, and they must be interpreted in terms of the standardized variables.

9.3.1 Identifying the Canonical Variables

Even though the canonical variables are artificial, they can be “identified” in terms of the subject-matter variables. Many times this identification is aided by computing the correlations between the canonical variates and the original variables. These correlations, however, must be interpreted with caution. They provide only univariate information, in the sense that they do not indicate how the original variables contribute jointly to the canonical analyses. For this reason, many investigators prefer to assess the contributions of the original variables directly from the standardized coefficients (9.2).

Let $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_p)'$ and $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_q)'$, so that the vectors of canonical variables are

$$\begin{array}{ccccccc} \mathbf{U} & = & \mathbf{A} & \mathbf{X}^{(1)} & \text{and} & \mathbf{V} & = & \mathbf{B} & \mathbf{X}^{(2)} \\ (p \times 1) & & (p \times p) & (p \times 1) & & (q \times 1) & & (q \times q) & (q \times 1) \end{array}$$

where we are primarily interested in the first p canonical variables in \mathbf{V} .

Introducing the $(p \times p)$ diagonal matrix $\mathbf{V}_{11}^{-1/2}$ with k th diagonal element

$$\sigma_{kk}^{-1/2} = \left(\text{Var}(X_k^{(1)}) \right)^{-1/2},$$

we find

$$\begin{aligned} (p \times p) \quad \boldsymbol{\rho}_{\mathbf{U}, \mathbf{X}^{(1)}} &= \text{Cor}(\mathbf{U}, \mathbf{X}^{(1)}) = \text{Cov}(\mathbf{U}, \mathbf{V}_{11}^{-1/2} \mathbf{X}^{(1)}) \\ &= \text{Cov}(\mathbf{A} \mathbf{X}^{(1)}, \mathbf{V}_{11}^{-1/2} \mathbf{X}^{(1)}) = \mathbf{A} \boldsymbol{\Sigma}_{11} \mathbf{V}_{11}^{-1/2}. \end{aligned} \quad (9.3)$$

Similar calculations for the pairs $(\mathbf{U}, \mathbf{X}^{(2)})$, $(\mathbf{V}, \mathbf{X}^{(2)})$ and $(\mathbf{V}, \mathbf{X}^{(1)})$ yield

$$\begin{aligned} (p \times q) \quad \boldsymbol{\rho}_{\mathbf{U}, \mathbf{X}^{(2)}} &= \mathbf{A} \boldsymbol{\Sigma}_{12} \mathbf{V}_{22}^{-1/2} \\ (q \times q) \quad \boldsymbol{\rho}_{\mathbf{V}, \mathbf{X}^{(2)}} &= \mathbf{B} \boldsymbol{\Sigma}_{22} \mathbf{V}_{22}^{-1/2} \\ (q \times p) \quad \boldsymbol{\rho}_{\mathbf{V}, \mathbf{X}^{(1)}} &= \mathbf{B} \boldsymbol{\Sigma}_{21} \mathbf{V}_{11}^{-1/2}, \end{aligned} \quad (9.4)$$

where $\mathbf{V}_{22}^{-1/2}$ is the $(q \times q)$ diagonal matrix with the i th diagonal element

$$\sigma_{ii}^{-1/2} = \left(\text{Var}(X_i^{(2)}) \right)^{-1/2}.$$

Canonical variables derived from standardized variables are sometimes interpreted by computing the correlations. Thus

$$\begin{aligned} \boldsymbol{\rho}_{\mathbf{U}, \mathbf{Z}^{(1)}} &= \mathbf{A}_Z \boldsymbol{\rho}_{11}, & \boldsymbol{\rho}_{\mathbf{V}, \mathbf{Z}^{(2)}} &= \mathbf{B}_Z \boldsymbol{\rho}_{22} \\ \boldsymbol{\rho}_{\mathbf{U}, \mathbf{Z}^{(2)}} &= \mathbf{A}_Z \boldsymbol{\rho}_{12}, & \boldsymbol{\rho}_{\mathbf{V}, \mathbf{Z}^{(1)}} &= \mathbf{B}_Z \boldsymbol{\rho}_{21} \end{aligned}$$

where \mathbf{A}_Z and \mathbf{B}_Z are the matrices whose rows contain the canonical coefficients for the $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)}$ sets, respectively.

Remark. The correlations are unaffected by the standardization, since for example

$$\boldsymbol{\rho}_{\mathbf{U}, \mathbf{X}^{(1)}} = \mathbf{A} \boldsymbol{\Sigma}_{11} \mathbf{V}_{11}^{-1/2} = \underbrace{\mathbf{A} \mathbf{V}_{11}^{1/2}}_{\mathbf{A}_Z} \underbrace{\mathbf{V}_{11}^{-1/2} \boldsymbol{\Sigma}_{11} \mathbf{V}_{11}^{-1/2}}_{\boldsymbol{\rho}_{11}} = \boldsymbol{\rho}_{\mathbf{U}, \mathbf{Z}^{(1)}}.$$

Remark. The correlations $\boldsymbol{\rho}_{U, \mathbf{X}^{(1)}}$ and $\boldsymbol{\rho}_{V, \mathbf{X}^{(2)}}$ can help supply meanings for the canonical variates. The spirit is the same as in principal component analysis when the correlations between the principal components and their associated variables may provide subject-matter interpretations for the components.

9.4 Sample Canonical Variates and Sample Canonical Correlations

A random sample of n observations on each of the $(p + q)$ variables $\mathbf{X}^{(1)}$, $\mathbf{X}^{(2)}$ can be assembled into the $n \times (p + q)$ data matrix

$$\mathbf{X} = (\mathbf{X}^{(1)} \mid \mathbf{X}^{(2)})$$

$$= \begin{pmatrix} x_{11}^{(1)} & \cdots & x_{1p}^{(1)} & \vdots & x_{11}^{(2)} & \cdots & x_{1q}^{(2)} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1}^{(1)} & \cdots & x_{np}^{(1)} & \vdots & x_{n1}^{(2)} & \cdots & x_{nq}^{(2)} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^{(1)'} & \vdots & \mathbf{x}_1^{(2)'} \\ \vdots & \vdots & \vdots \\ \mathbf{x}_n^{(1)'} & \vdots & \mathbf{x}_n^{(2)'} \end{pmatrix}.$$

We find $\bar{\mathbf{x}} = \begin{pmatrix} \bar{\mathbf{x}}^{(1)} \\ \bar{\mathbf{x}}^{(2)} \end{pmatrix}$ where $\bar{\mathbf{x}}^{(i)} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^{(i)}$, $i = 1, 2$ and

$$\mathbf{S} = \begin{pmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{pmatrix} \begin{matrix} (p \times (p+q)) \\ (q \times (p+q)) \end{matrix} \quad \text{with} \quad \mathbf{S}_{12} = \mathbf{S}_{21}'$$

$$\begin{matrix} ((p+q) \times p) & ((p+q) \times q) \end{matrix}$$

and

$$\mathbf{S}_{kl} = \frac{1}{n-1} \sum_{j=1}^n (\mathbf{x}_j^{(k)} - \bar{\mathbf{x}}^{(k)})(\mathbf{x}_j^{(l)} - \bar{\mathbf{x}}^{(l)})', \quad k, l = 1, 2. \quad (9.5)$$

The linear combinations $\hat{U} = \hat{\mathbf{a}}' \mathbf{x}^{(1)}$ and $\hat{V} = \hat{\mathbf{b}}' \mathbf{x}^{(2)}$ have sample correlation

$$r_{\hat{U}, \hat{V}} = \frac{\hat{\mathbf{a}}' \mathbf{S}_{12} \hat{\mathbf{b}}}{\sqrt{\hat{\mathbf{a}}' \mathbf{S}_{11} \hat{\mathbf{a}}} \sqrt{\hat{\mathbf{b}}' \mathbf{S}_{22} \hat{\mathbf{b}}}}. \quad (9.6)$$

The first pair of sample canonical variates is the pair of linear combinations \hat{U}_1 , \hat{V}_1 having unit sample variances that maximizes the ratio (9.6).

In general, the k th pair of sample canonical variates is the pair of linear combinations \hat{U}_k , \hat{V}_k having unit sample variances that maximizes the ratio (9.6) among those linear combinations uncorrelated with the previous $k - 1$ sample canonical variates.

The sample correlation between \hat{U}_k and \hat{V}_k is called the k th sample canonical correlation.

Example (Soil Evaporation, p. 1-8). Table 9.1 and Figure 9.1 show the results of the canonical correlation analysis of the soil and the air variables.

Table 9.1: Results of the canonical correlation analysis. Data set: Soil Evaporation, p. 1-8.

```
> cc.airsoil
$cor:
[1] 0.9624326 0.7604630 0.5963187

$xcoef:
      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
[1,] -0.00810700994 -0.0131838494 -0.0161725757  0.02517594372  0.04412103124 -0.0579086504 -0.0005754910
[2,]  0.00216809356 -0.0440762196 -0.0176049769  0.01160575703 -0.08327800431 -0.0336352188 -0.0366291212
[3,] -0.00470274969  0.0086711986  0.0003125829 -0.00638888876  0.00565562936  0.0235464153  0.0041274285
[4,]  0.01681633315 -0.0389196441  0.0407784120  0.12108216355 -0.05106799333 -0.0318368973  0.0450208056
[5,]  0.01080944196 -0.0294491540  0.0218271660 -0.00541588722  0.02269311835 -0.0180959750  0.0190990783
[6,] -0.00218121839  0.0096778587 -0.0110017742  0.00148441044  0.00070780115  0.0098347719 -0.0065376621
[7,]  0.00001586807  0.0005137182 -0.0004222563  0.00007431056  0.00004439516 -0.0004388761  0.0009371029

$ycoef:
      [,1]      [,2]      [,3]
[1,] -0.023889982  0.079116806 -0.02793426
[2,]  0.006589414  0.007931868 -0.13448882
[3,] -0.001183075 -0.026269281  0.02816080

$xcen:
      maxat  minat  avat      maxh  minh      avh      wind
90.73913  70.06522  190.5  94.71739  48.5  396.913  277.6739

$ycen:
      maxst  minst  avst
87.56522  71.26087  173.5217
```

Proposition 9.4.1. Let $\hat{\rho}_1^* \geq \dots \geq \hat{\rho}_p^*$ be the p ordered eigenvalues of

$$\mathbf{S}_{11}^{-1/2} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1/2}$$

with corresponding eigenvectors $\hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_p$, where the \mathbf{S}_{kl} are defined in (9.5) and $p \leq q$. Let $\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_p$ be the eigenvectors of

$$\mathbf{S}_{22}^{-1/2} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12} \mathbf{S}_{22}^{-1/2},$$

where the first p eigenvectors $\hat{\mathbf{f}}$ may be obtained from

$$\hat{\mathbf{f}}_k = \frac{1}{\hat{\rho}_k^*} \mathbf{S}_{22}^{-1/2} \mathbf{S}_{21} \mathbf{S}_{11}^{-1/2} \hat{\mathbf{e}}_k, \quad k = 1, \dots, p.$$

Then the k th sample canonical variate pair is

$$\hat{U}_k = \underbrace{\hat{\mathbf{e}}_k' \mathbf{S}_{11}^{-1/2}}_{=\hat{\mathbf{a}}_k'} \mathbf{x}^{(1)} \quad \hat{V}_k = \underbrace{\hat{\mathbf{f}}_k' \mathbf{S}_{22}^{-1/2}}_{=\hat{\mathbf{b}}_k'} \mathbf{x}^{(2)}$$

where $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ are the values of the variables $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ for a particular experimental unit.

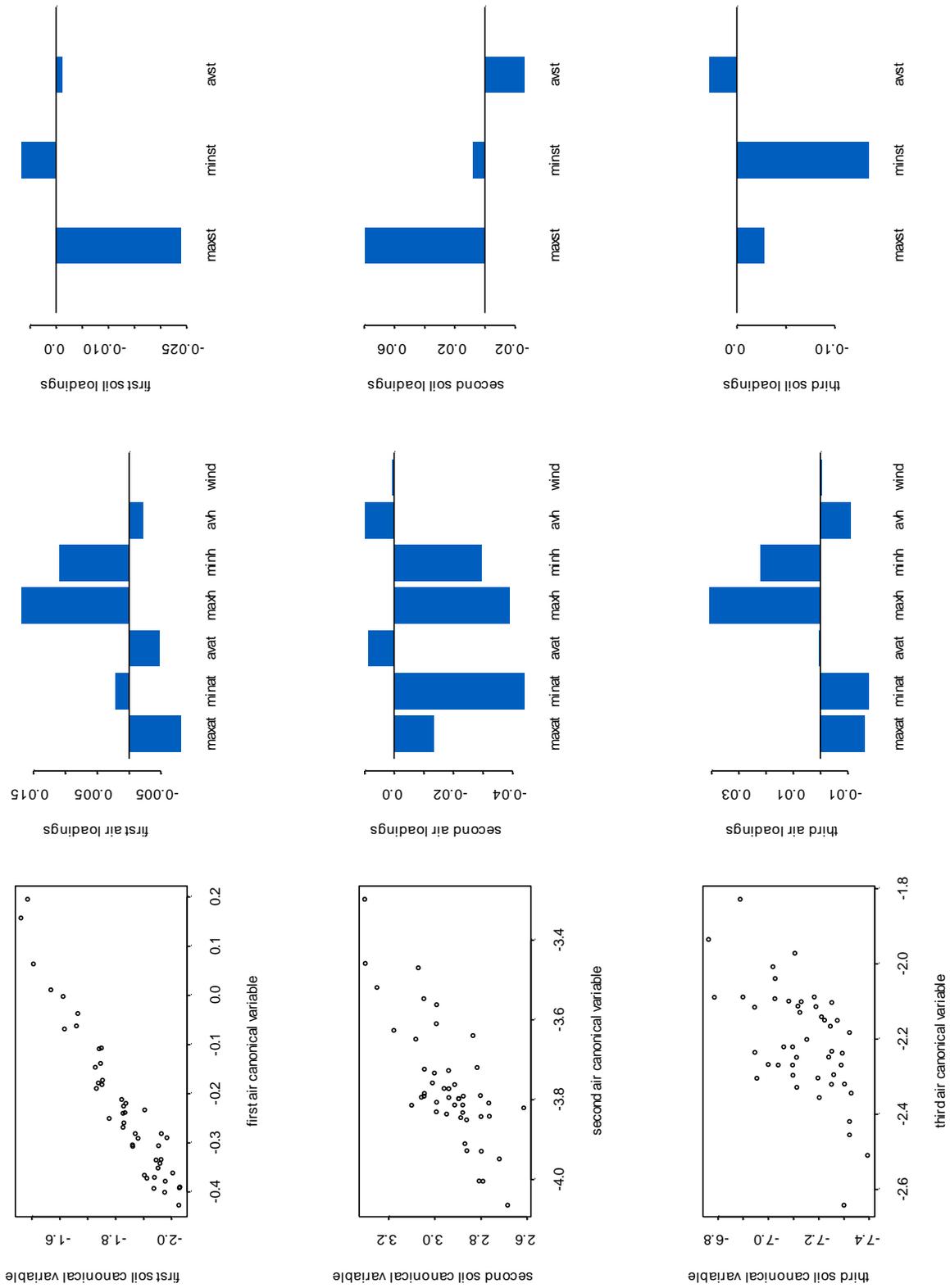


Figure 9.1: Results of the canonical correlation analysis. Data set: Soil Evaporation, p. 1-8.

Remark. The interpretation of \hat{U}_k and \hat{V}_k is often aided by computing the sample correlations between the canonical variates and the variables in the sets $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$. We define the matrices

$$\hat{\mathbf{A}} = \begin{matrix} & (\hat{\mathbf{a}}_1 & , \dots , & \hat{\mathbf{a}}_p)' \\ (p \times p) & (p \times 1) & & (p \times 1) \end{matrix}$$

$$\hat{\mathbf{B}} = \begin{matrix} & (\hat{\mathbf{b}}_1 & , \dots , & \hat{\mathbf{b}}_q)' \\ (q \times q) & (q \times 1) & & (q \times 1) \end{matrix}$$

whose rows are the coefficient vectors for the sample canonical variates. We find

$$\hat{\mathbf{U}} = \begin{matrix} & \hat{\mathbf{A}} & \mathbf{x}^{(1)} \\ (p \times 1) & (p \times p) & (p \times 1) \end{matrix} \quad \text{and} \quad \hat{\mathbf{V}} = \begin{matrix} & \hat{\mathbf{B}} & \mathbf{x}^{(2)} \\ (q \times 1) & (q \times q) & (q \times 1) \end{matrix}$$

and can define

$$\begin{aligned} \mathbf{R}_{\hat{\mathbf{U}}, \mathbf{x}^{(1)}} &= \text{matrix of sample correlations of } \hat{\mathbf{U}} \text{ with } \mathbf{x}^{(1)} \\ \mathbf{R}_{\hat{\mathbf{V}}, \mathbf{x}^{(2)}} &= \text{matrix of sample correlations of } \hat{\mathbf{V}} \text{ with } \mathbf{x}^{(2)} \\ \mathbf{R}_{\hat{\mathbf{U}}, \mathbf{x}^{(2)}} &= \text{matrix of sample correlations of } \hat{\mathbf{U}} \text{ with } \mathbf{x}^{(2)} \\ \mathbf{R}_{\hat{\mathbf{V}}, \mathbf{x}^{(1)}} &= \text{matrix of sample correlations of } \hat{\mathbf{V}} \text{ with } \mathbf{x}^{(1)} \end{aligned}$$

and corresponding to (9.3) and (9.4) we find

$$\begin{aligned} \mathbf{R}_{\hat{\mathbf{U}}, \mathbf{x}^{(1)}} &= \hat{\mathbf{A}} \mathbf{S}_{11} \mathbf{D}_{11}^{-1/2} \\ \mathbf{R}_{\hat{\mathbf{V}}, \mathbf{x}^{(2)}} &= \hat{\mathbf{B}} \mathbf{S}_{22} \mathbf{D}_{22}^{-1/2} \\ \mathbf{R}_{\hat{\mathbf{U}}, \mathbf{x}^{(2)}} &= \hat{\mathbf{A}} \mathbf{S}_{12} \mathbf{D}_{22}^{-1/2} \\ \mathbf{R}_{\hat{\mathbf{V}}, \mathbf{x}^{(1)}} &= \hat{\mathbf{B}} \mathbf{S}_{21} \mathbf{D}_{11}^{-1/2}, \end{aligned}$$

where $\mathbf{D}_{11}^{-1/2}$ and $\mathbf{D}_{22}^{-1/2}$ are the $(p \times p)$ and $(q \times q)$ diagonal matrix with i th diagonal element

$$\left(\text{sample variance}(x_i^{(1)}) \right)^{-1/2} \quad \text{and} \quad \left(\text{sample variance}(x_i^{(2)}) \right)^{-1/2},$$

respectively.

9.5 Canonical Correlation Analysis applied to Fields and Forecasting with Canonical Correlation Analysis

Source: Wilks (2006) pp. 519-522 (see next pages).

12.2 CCA Applied to Fields

12.2.1 Translating Canonical Vectors to Maps

Canonical correlation analysis is usually most interesting for atmospheric data when applied to fields. Here the spatially distributed observations (either at gridpoints or observing locations) are encoded into the vectors \mathbf{x} and \mathbf{y} in the same way as for PCA. That is, even though the data may pertain to a two- or three-dimensional field, each location is numbered sequentially and pertains to one element of the corresponding data vector. It is not necessary for the spatial domains encoded into \mathbf{x} and \mathbf{y} to be the same, and indeed in the applications of CCA that have appeared in the literature they are usually different.

As is the case with the use of PCA with spatial data, it is often informative to plot maps of the canonical vectors by associating the magnitudes of their elements and the geographic locations to which they pertain. In this context the canonical vectors are sometimes called canonical patterns, since the resulting maps show spatial patterns of the ways in which the original variables contribute to the canonical variables. Examining the pairs of maps formed by corresponding vectors \mathbf{a}_m and \mathbf{b}_m can be informative about the nature of the relationship between variations in the data over the two domains encoded in \mathbf{x} and \mathbf{y} , respectively. Figures 12.2 and 12.3 show examples of maps of canonical vectors.

It can also be informative to plot pairs of maps of the homogeneous (Equation 12.7) or heterogeneous correlations (Equation 12.8). Each of these vectors contain correlations between an underlying data field and one of the canonical variables, and these correlations can also be plotted at the corresponding locations. Figure 12.1, from Wallace *et al.* (1992), shows one such pair of homogeneous correlation patterns. Figure 12.1a shows the spatial distribution of correlations between a canonical variable v , and the values of the corresponding data \mathbf{x} that contains values of average December-February sea-surface temperatures (SSTs) in the north Pacific Ocean. This canonical variable accounts for 18% of the total variance of the SSTs in the data set analyzed (Equation 12.12). Figure 12.1b shows the spatial distribution of the correlations for the corresponding canonical variable w , that pertains to average hemispheric 500 mb heights \mathbf{y} during the same winters included in the SST data in \mathbf{x} . This canonical variable accounts for 23% of the total variance of the winter hemispheric height variations. The correlation pattern in Figure 12.1a corresponds to either cold water in the central north Pacific and warm water along the west coast of North America, or warm water in the central north Pacific and cold water along the west coast of North America. The pattern of 500 mb height correlations in Figure 12.1b is remarkably similar to the PNA pattern (cf. Figures 11.10b and 3.28).

The correlation between the two time series v and w is the canonical correlation $r_C = 0.79$. Because v and w are well correlated, these figures indicate that cold SSTs in the central Pacific simultaneously with warm SSTs in the northeast Pacific (relatively large positive v) tend to coincide with a 500 mb ridge over northwestern North America and a 500 mb trough over southeastern North America (relatively large positive w). Similarly, warm water in the central north Pacific and cold water in the northwestern Pacific (relatively large negative v) are associated with the more zonal PNA flow (relatively large negative w).

12.2.2 Combining CCA with PCA

The sampling properties of CCA can be poor when the available data are few relative to the dimensionality of the data vectors. The result can be that sample estimates for

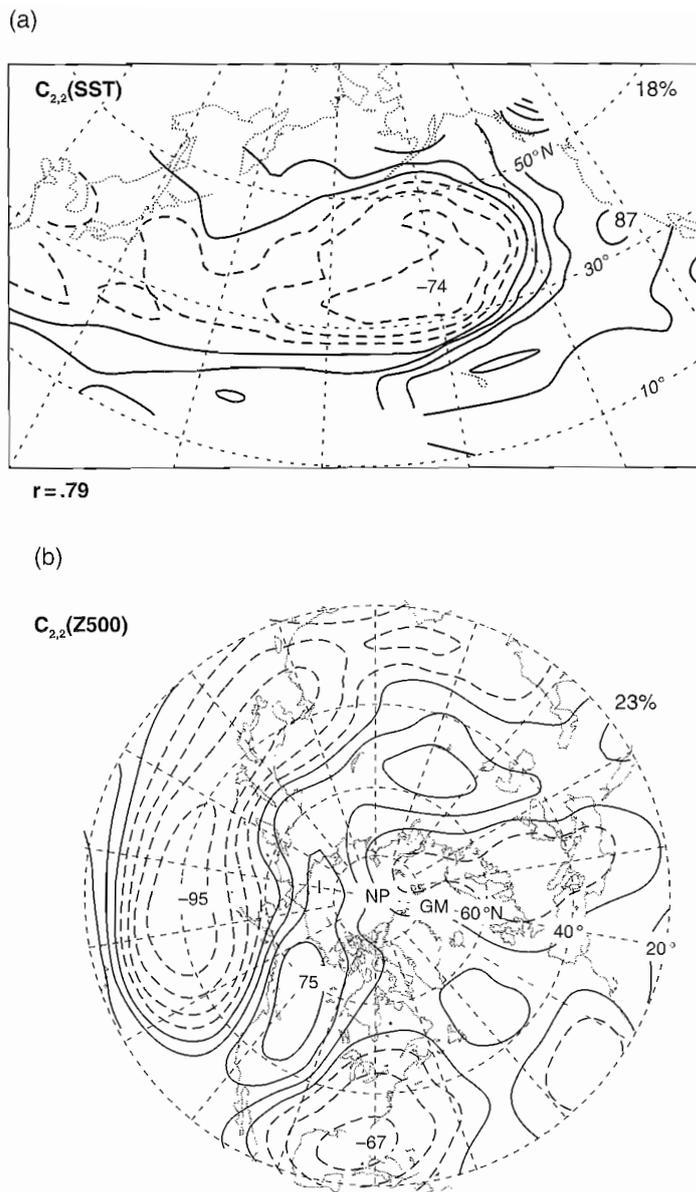


FIGURE 12.1 Homogeneous correlation maps for a pair of canonical variables pertaining to (a) average winter sea-surface temperatures (SSTs) in the northern Pacific Ocean, and (b) hemispheric winter 500 mb heights. The pattern of SST correlation in the left-hand panel (and its negative) are associated with the PNA pattern of 500 mb height correlations shown in the right-hand panel. The canonical correlation for this pair of canonical variables is 0.79. From Wallace *et al.* (1992).

CCA parameters may be unstable (i.e., exhibit large variations from batch to batch) for small samples (e.g., Bretherton *et al.* 1992; Cherry 1996; Friederichs and Hense 2003). Friederichs and Hense (2003) describe, in the context of atmospheric data, both conventional parametric tests and resampling tests to help assess whether sample canonical correlations may be spurious sampling artifacts. These tests examine the null hypothesis that all the underlying population canonical correlations are zero.

Relatively small sample sizes are common when analyzing time series of atmospheric fields. In CCA, it is not uncommon for there to be fewer observations n than the dimensions I and J of the data vectors, in which case the necessary matrix inversions

cannot be computed (see Section 12.3). However, even if the sample sizes are large enough to carry through the calculations, sample CCA statistics are erratic unless $n \gg M$. Barnett and Preisendorfer (1987) suggested that a remedy for this problem is to prefilter the two fields of raw data using separate PCAs before subjecting them to a CCA, and this has become a conventional procedure. Rather than directly correlating linear combinations of the fields \mathbf{x}' and \mathbf{y}' , the CCA operates on the vectors \mathbf{u}_x and \mathbf{u}_y , which consist of the leading principal components of \mathbf{x}' and \mathbf{y}' . The truncations for these two PCAs (i.e., the dimensions of the vectors \mathbf{u}_x and \mathbf{u}_y) need not be the same, but should be severe enough for the larger of the two to be substantially smaller than the sample size n . Livezey and Smith (1999) provide some guidance for the subjective choices that need to be made in this approach.

This combined PCA/CCA approach is not always best, and can be inferior if important information is discarded when truncating the PCA. In particular, there is no guarantee that the most strongly correlated linear combinations of \mathbf{x} and \mathbf{y} will be well related to the leading principal components of one field or the other.

10 Discrimination and Classification

10.1 Introduction

Discrimination and classification are multivariate techniques concerned with separating distinct sets of objects (or observations) and with allocating new objects (observations) to previously defined groups. Discriminant analysis is rather exploratory in nature. As a separative procedure, it is often employed on a one-time basis in order to investigate observed differences when causal relationships are not well understood. Classification procedures are less exploratory in the sense that they lead to well-defined rules, which can be used for assigning new objects. Classification ordinarily requires more problem structure than discrimination does.

Thus, the immediate goals of discrimination and classification, respectively, are as follows:

- Discrimination (or separation): describes the differential features of objects (observations) from several known collections (populations). We try to find “discriminants” whose numerical values are such that the collection is separated as much as possible.
- Classification (or allocation): sorts objects (observations) into two or more labeled classes. The emphasis is on deriving a rule that can be used to optimally assign new objects to the labeled classes.

In practice, the two goals frequently overlap, and the distinction between separation and allocation becomes blurred.

10.2 Separation and Classification for Two Populations

Example (El Niño). Consider the data set in Table 10.1. With discriminant analysis we want to classify the years in the data set as either El Niño or non-El Niño, on the basis of the corresponding temperature and pressure data.

It is convenient to label the classes π_1 and π_2 . The objects are ordinarily separated or classified on the basis of measurements on, for instance, p associated random variables $\mathbf{X}' = (X_1, \dots, X_p)$. The observed values of \mathbf{X} differ to some extent from one class to the other. We can think of the totality of values from the first class as being the population \mathbf{x} values for π_1 and those from the second class as the population of \mathbf{x} values for π_2 . These two populations can then be described by probability density functions $f_1(\mathbf{x})$ and

Table 10.1: June climate data for Guayaquil, Ecuador, 1951-1970. Source: Wilks (2006).

Year	Temp. (°C)	Prec. (mm)	Pressure (hPa)	El Nino yes=1, no=0
1951	26.1	43	1009.5	1
1952	24.5	10	1010.9	0
1953	24.8	4	1010.7	1
1954	24.5	0	1011.2	0
1955	24.1	2	1011.9	0
1956	24.3	NA	1011.2	0
1957	26.4	31	1009.3	1
1958	24.9	0	1011.1	0
1959	23.7	0	1012.0	0
1960	23.5	0	1011.4	0
1961	24.0	2	1010.9	0
1962	24.1	3	1011.5	0
1963	23.7	0	1011.0	0
1964	24.3	4	1011.2	0
1965	26.6	15	1009.9	1
1966	24.6	2	1012.5	0
1967	24.8	0	1011.1	0
1968	24.4	1	1011.8	0
1969	26.8	127	1009.3	1
1970	25.2	2	1010.6	0

$f_2(\mathbf{x})$, and consequently, we can talk of assigning observations to populations or objects to classes interchangeably.

A good classification procedure should:

- result in a few misclassifications. In other words, the chances, or probabilities, of misclassification should be small
- take into account that one class or population has a greater likelihood of occurrence than another because one of the two populations is relatively much larger than the other

Example. A randomly selected firm should be classified as nonbankrupt unless the data overwhelmingly favors bankruptcy.

- whenever possible, account for the costs associated with misclassification.

Classification rules cannot usually provide an error-free method of assignment. This is because there may not be a clear distinction between the measured characteristics of

the populations; that is, the groups may overlap. It is then possible, for example, to incorrectly classify a π_2 object as belonging to π_1 or a π_1 object as belonging to π_2 .

Let $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ be the probability density functions associated with the $(p \times 1)$ random vector \mathbf{X} for the populations π_1 and π_2 , respectively. An object with associated measurements \mathbf{x} must be assigned to either π_1 or π_2 . Let Ω be the sample space and R_1 that set of \mathbf{x} values for which we classify objects as π_1 and $R_2 = \Omega - R_1$ be the remaining \mathbf{x} values for which we classify objects as π_2 . Since every object must be assigned to one and only one of the two populations, the sets R_1 and R_2 are mutually exclusive and exhaustive. For $p = 2$, we might have a case like the one pictured in Figure 10.1.

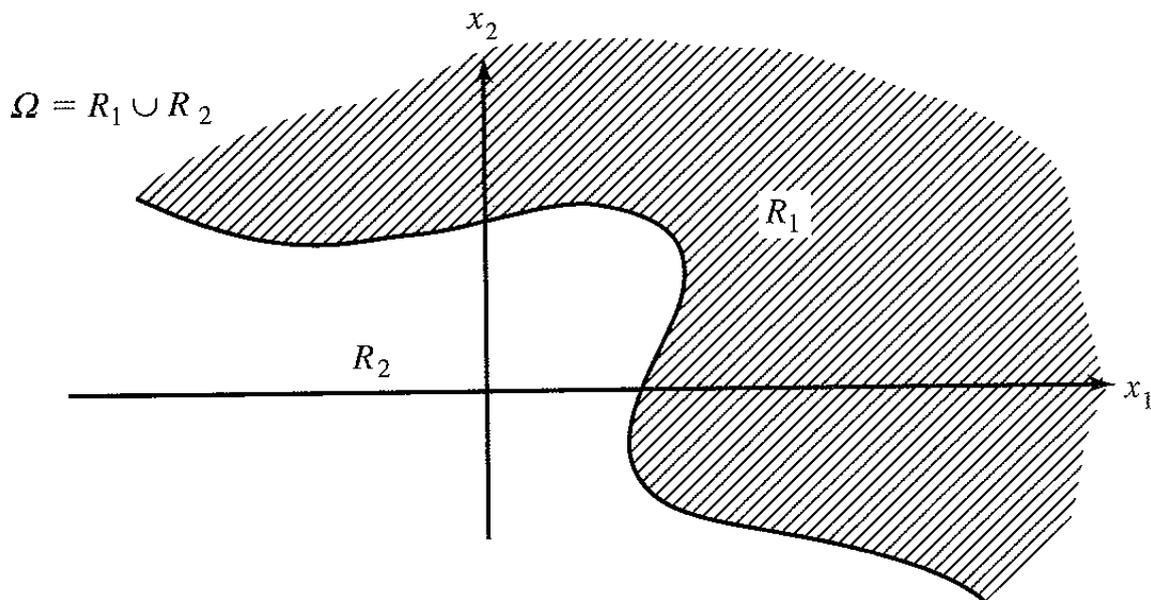


Figure 10.1: Classification regions for two populations. Source: Johnson and Wichern (2007).

The conditional probability, $P(2|1)$, of classifying an object as π_2 when, in fact, it is from π_1 is

$$P(2|1) = P(\mathbf{X} \in R_2 | \pi_1) = \int_{R_2} f_1(\mathbf{x}) d\mathbf{x}.$$

Similarly, the conditional probability, $P(1|2)$, of classifying an object as π_1 when it is really from π_2 is

$$P(1|2) = P(\mathbf{X} \in R_1 | \pi_2) = \int_{R_1} f_2(\mathbf{x}) d\mathbf{x}.$$

This is illustrated in Figure 10.2 for the univariate case, $p = 1$.

Let $p_1 = P(\pi_1)$ be the prior probability of π_1 and $p_2 = P(\pi_2)$ be the prior probability of π_2 , where $p_1 + p_2 = 1$. Then the overall probabilities of correctly or incorrectly

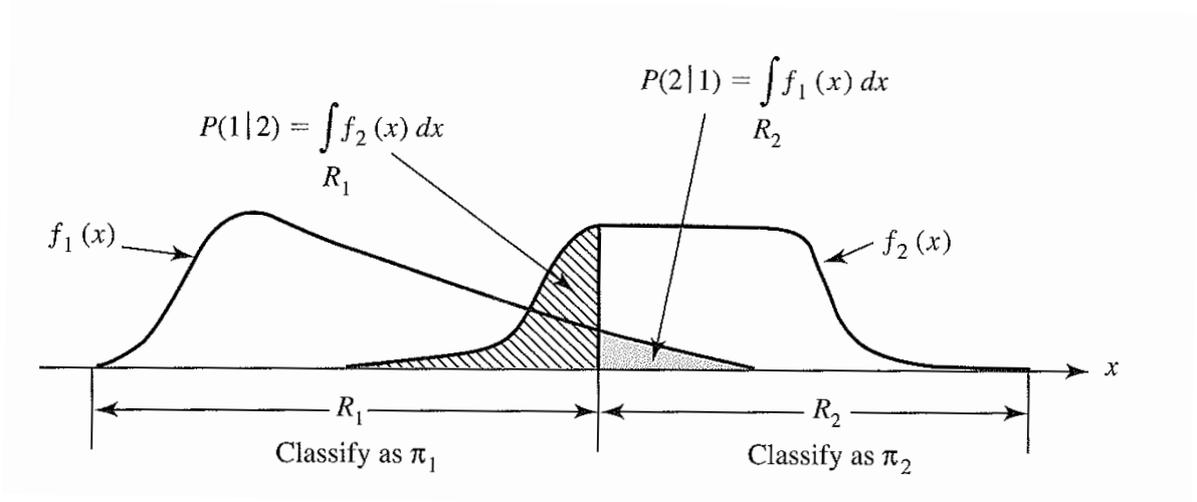


Figure 10.2: Misclassification probabilities for hypothetical classification regions. Source: Johnson and Wichern (2007).

classifying objects can be derived as the product of the prior and conditional classification probabilities:

$$\begin{aligned}
 &P(\text{observation is correctly classified as } \pi_1) \\
 &= P(\text{observation comes from } \pi_1 \text{ and is correctly classified as } \pi_1) \\
 &= P(\mathbf{X} \in R_1 | \pi_1)P(\pi_1) = P(1|1)p_1 \\
 &P(\text{observation is misclassified as } \pi_1) \\
 &= P(\text{observation comes from } \pi_2 \text{ and is misclassified as } \pi_1) \\
 &= P(\mathbf{X} \in R_1 | \pi_2)P(\pi_2) = P(1|2)p_2 \\
 &P(\text{observation is correctly classified as } \pi_2) \\
 &= P(\text{observation comes from } \pi_2 \text{ and is correctly classified as } \pi_2) \\
 &= P(\mathbf{X} \in R_2 | \pi_2)P(\pi_2) = P(2|2)p_2 \\
 &P(\text{observation is misclassified as } \pi_2) \\
 &= P(\text{observation comes from } \pi_1 \text{ and is misclassified as } \pi_2) \\
 &= P(\mathbf{X} \in R_2 | \pi_1)P(\pi_1) = P(2|1)p_1.
 \end{aligned}$$

Classification schemes are often evaluated in terms of their misclassification probabilities, but this ignores misclassification costs. A rule that ignores costs may cause problems. The costs of misclassification can be defined by a cost matrix (Table 10.2)

The costs are zero for correct classification, $c(1|2)$ when an observation from π_2 is incorrectly classified as π_1 , and $c(2|1)$ when a π_1 observation is incorrectly classified as π_2 .

For any rule, the average, or expected cost of misclassification (ECM) is provided by multiplying the off-diagonal entries in Table 10.2 by their probabilities of occurrence.

Table 10.2: Costs of misclassification

		classify as:	
		π_1	π_2
true population:	π_1	0	$c(2 1)$
	π_2	$c(1 2)$	0

Consequently,

$$\text{ECM} = c(2|1)P(2|1)p_1 + c(1|2)P(1|2)p_2.$$

A reasonable classification rule should have an ECM as small as possible.

Proposition 10.2.1. *The regions R_1 and R_2 that minimize the ECM are defined by the values \mathbf{x} for which the following inequalities hold:*

$$\begin{aligned}
 R_1 : \quad & \underbrace{\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})}}_{\text{density ratio}} \geq \underbrace{\frac{c(1|2)}{c(2|1)}}_{\text{cost ratio}} \underbrace{\frac{p_2}{p_1}}_{\text{prior probability ratio}} \\
 R_2 : \quad & \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1}
 \end{aligned} \tag{10.1}$$

Remark. It is interesting to consider the classification regions defined in (10.1) for some special cases:

- When the prior probabilities are unknown, they are often taken to be equal, i.e. $p_2/p_1 = 1$, and the minimum ECM rule involves comparing the ratio of the population densities to the ratio of the appropriate misclassification costs;
- equal misclassification costs: $c(1|2)/c(2|1) = 1$;
- equal prior probabilities and equal misclassification costs: $p_2/p_1 = c(1|2)/c(2|1) = 1$.

If \mathbf{x}_0 is a new observation, it is assigned to π_1 , if

$$\frac{f_1(\mathbf{x}_0)}{f_2(\mathbf{x}_0)} \geq \frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1}.$$

On the other hand, if

$$\frac{f_1(\mathbf{x}_0)}{f_2(\mathbf{x}_0)} < \frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1},$$

we assign \mathbf{x}_0 to π_2 .

Remark. There are other criterias to derive “optimal” classification procedures. For example, choose R_1 and R_2 to minimize the total probability of misclassification (TPM):

$$\begin{aligned} \text{TPM} &= P(\text{misclassifying a } \pi_1 \text{ observation or misclassifying a } \pi_2 \text{ observation}) \\ &= P(\text{observation comes from } \pi_1 \text{ and is misclassified}) \\ &\quad + P(\text{observation comes from } \pi_2 \text{ and is misclassified}) \\ &= p_1 \int_{R_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

Mathematically, this problem is equivalent to minimizing the expected cost of misclassification when the costs of misclassification are equal.

10.2.1 Classification with Two Multivariate Normal Populations

Classification procedures based on normal populations predominate in statistical practice because of their simplicity and reasonably high efficiency across a wide variety of population models. We now assume that $f_1(\mathbf{x})$ and $f_2(\mathbf{x})$ are multivariate normal densities, the first with mean vector $\boldsymbol{\mu}_1$ and covariance matrix $\boldsymbol{\Sigma}_1$ and the second with mean vector $\boldsymbol{\mu}_2$ and covariance matrix $\boldsymbol{\Sigma}_2$. There are two cases to distinguish:

i) Case $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}$

Suppose that the joint densities of $\mathbf{X}' = (X_1, \dots, X_p)$ for population π_1 and π_2 are given by

$$f_i(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-1/2(\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right) \quad \text{for } i = 1, 2. \quad (10.2)$$

- Suppose that $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$, and $\boldsymbol{\Sigma}$ are known. Then, the minimum ECM regions in (10.1) become

$$\begin{aligned} R_1 : & \exp\left(-1/2(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + 1/2(\mathbf{x} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)\right) \\ & \geq \frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \\ R_2 : & \exp\left(-1/2(\mathbf{x} - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + 1/2(\mathbf{x} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_2)\right) \\ & < \frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \end{aligned}$$

Proposition 10.2.2. *Let the populations π_1 and π_2 be described by multivariate normal densities of the form (10.2). Then the allocation rule that minimizes the ECM is as follows: Allocate \mathbf{x}_0 to π_1 if*

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} \mathbf{x}_0 - 1/2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) \geq \ln\left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1}\right). \quad (10.3)$$

Allocate \mathbf{x}_0 to π_2 otherwise.

10.2.2 Fisher's Approach to Classification with Two Populations

Fisher's idea was to transform the multivariate observations \mathbf{x} to univariate observations y such that the y 's derived from population π_1 and π_2 were separated as much as possible. Fisher suggested taking linear combinations of \mathbf{x} to create y 's because they are simple enough functions of the \mathbf{x} to be handled easily. Fisher's approach does not assume that the populations are normal. It does, however, implicitly assume that the population covariance matrices are equal.

Proposition 10.2.4. *The linear combination $\hat{y} = \hat{\mathbf{a}}'\mathbf{x} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\mathbf{S}_{pooled}^{-1}\mathbf{x}$ maximizes the ratio*

$$\frac{\text{squared distance between sample means of } y}{\text{sample variance of } y} = \frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} = \frac{(\hat{\mathbf{a}}'\bar{\mathbf{x}}_1 - \hat{\mathbf{a}}'\bar{\mathbf{x}}_2)^2}{\hat{\mathbf{a}}'\mathbf{S}_{pooled}\hat{\mathbf{a}}}$$

over all possible coefficient vectors $\hat{\mathbf{a}}$ (see Figure 10.3). The maximum of the ratio is

$$D^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'\mathbf{S}_{pooled}^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2).$$

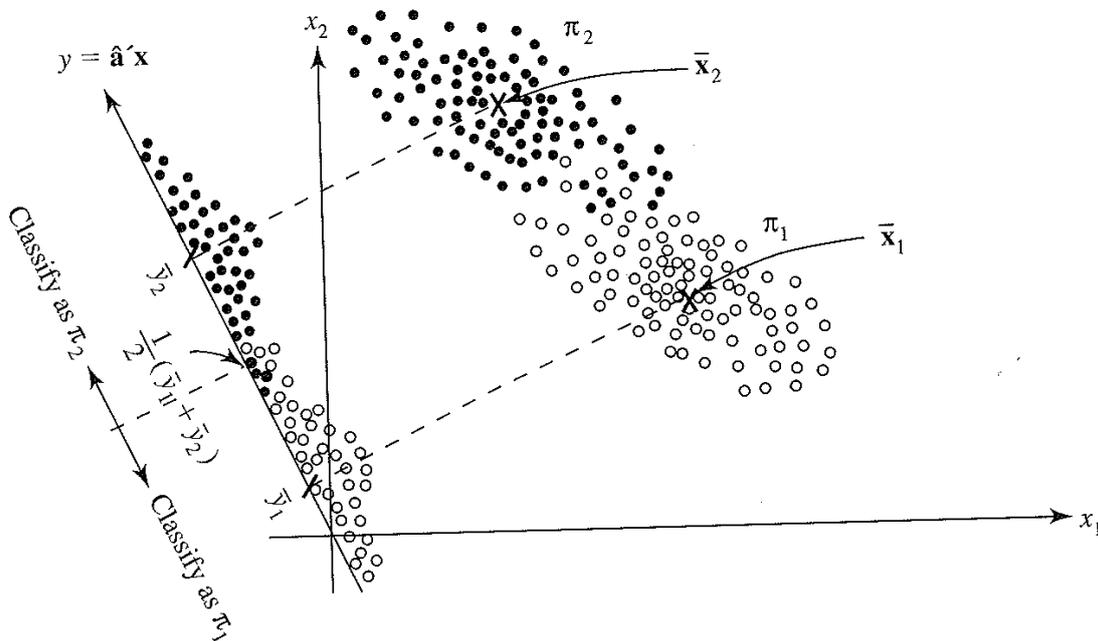


Figure 10.3: Fisher's procedure for two populations with $p = 2$. Source: Johnson and Wichern (2007).

Proposition 10.2.5. *An allocation rule based on Fisher's discriminant function: Allocate \mathbf{x}_0 to π_1 if*

$$\begin{aligned}\hat{y}_0 &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 \\ &\geq \hat{m} = 1/2(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2).\end{aligned}$$

Allocate \mathbf{x}_0 to π_2 if

$$\hat{y}_0 < \hat{m}.$$

Example (El Niño, p. 10-2). In Figure 10.4 the Fisher discriminant function is shown for the pairs of variables pressure (Pres), temperature (Temp) and precipitation (Prec), temperature (Temp), respectively. There is one misclassification, namely the year 1953.

Remark. The performance of any classification procedure should always be checked. Ideally, there will be enough data available to provide for “training” samples and “validation” samples. The training samples can be used to develop the classification function, and the validation samples can be used to evaluate its performance.

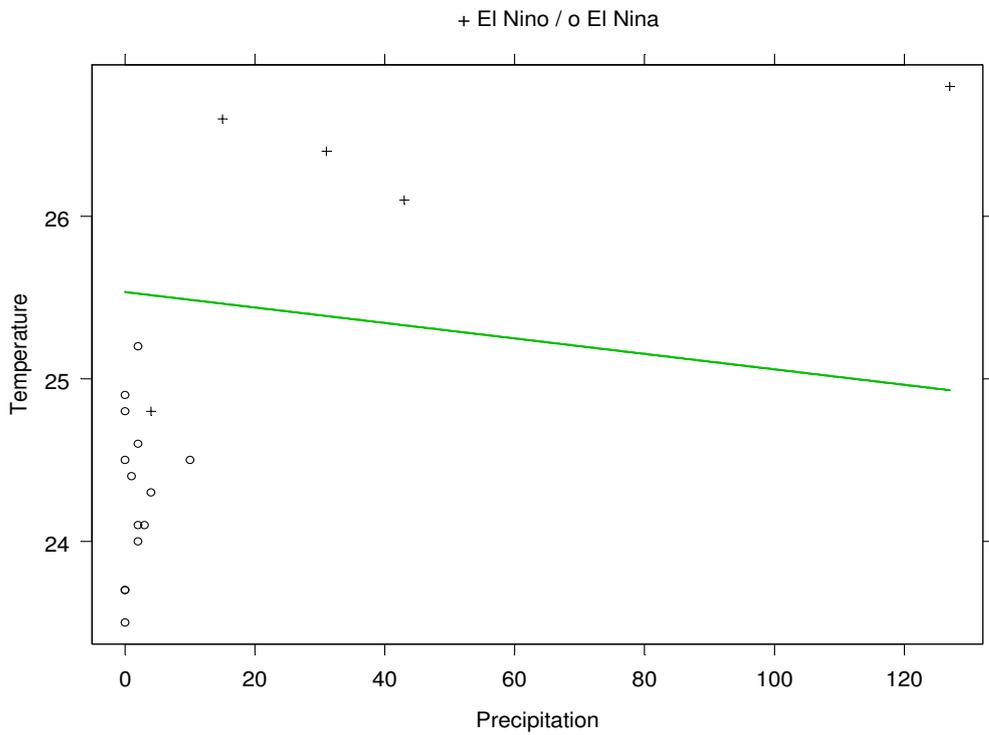
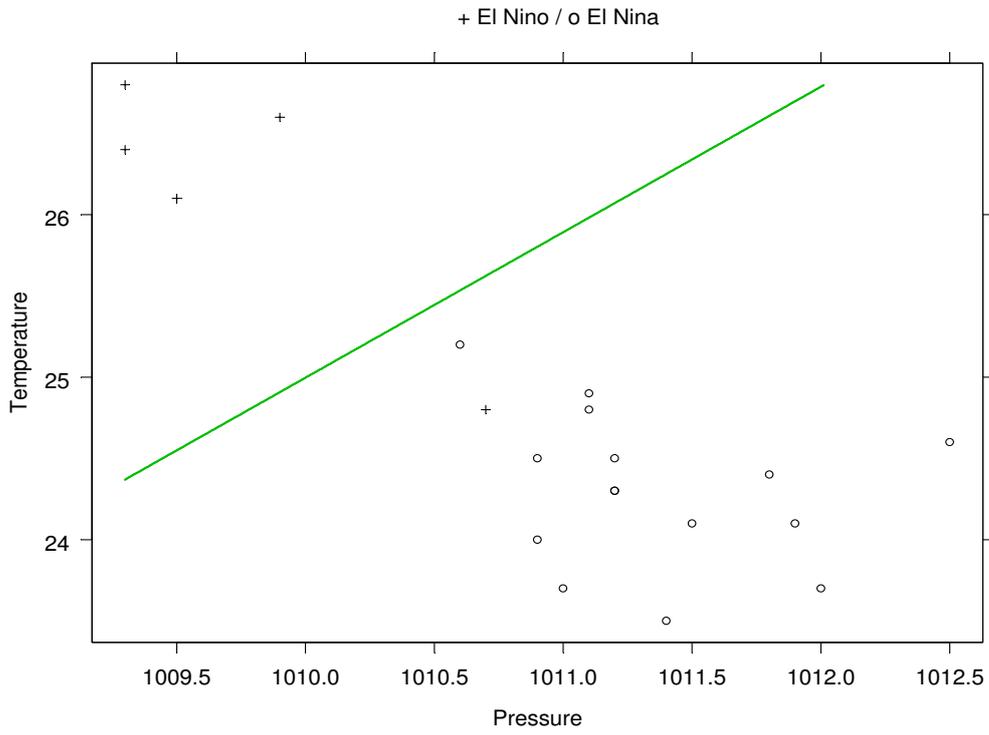


Figure 10.4: Discriminant function of the data for Guayaquil. Data set: Table 10.1, p. 10-2.