

RENATO LEONI

Principal Component Analysis

**UNIVERSITY OF FLORENCE
DEPARTMENT OF STATISTICS "G. PARENTI"
FLORENCE, 2007**

This paper is intended for a personal use only. Trading is not allowed.

1 INTRODUCTION

Researchers are frequently faced with the task of analysing a data collection concerning a large number of quantitative variables measured on many individuals (units) and usually displayed in tabular form. The aim of the analysis is often to find out patterns of interrelationships which may exist among variables or individuals. The problem is that, given the data volume, this aim is not readily achieved.

The focus of principal component analysis (PCA) is on the study of a large data collection of the type mentioned above from the point of view of the interrelationships which may subsist among variables or individuals, providing at the same time the researcher with a graphical representation of results on a subspace of low dimension (usually one or two).

In this paper, without making any assumption about an underlying probabilistic model, we will present the main features of PCA.

The contents of the paper can be summarized as follows.

In Section 2, the basic data and their algebraic structure are set out. In Section 3 and 4, privileging a geometrical language, some concepts which will be used extensively during the paper are introduced. Section 5 is devoted to a presentation of an approach to PCA. In Section 6, rules for a graphical representation of results are given. Finally, in Section 7, other approaches to PCA are set out ⁽¹⁾.

(1) Numerical examples, based both on fictitious and real data, are provided apart. Relevant algebraic concepts are stated in [20].

2 BASIC DATA AND THEIR ALGEBRAIC STRUCTURE

2.1 RAW DATA MATRIX

Consider the matrix (*raw data matrix*)

$$\mathbf{X} = \begin{bmatrix} X_{11} & \cdots & X_{1p} \\ \cdots & \cdots & \cdots \\ X_{n1} & \cdots & X_{np} \end{bmatrix}$$

where x_{ij} ($i = 1, \dots, n; j = 1, \dots, p$) denotes the value of the j th quantitative variable observed on the i th individual.

Although in practical applications the number n of individuals is often strictly greater of the number p of variables, that assumption is not necessary in performing PCA and will be dropped; in other words, we will suppose that it may be $n \gtrless p$.

Notice that, setting ($i = 1, \dots, n$)

$$\mathbf{x}_i = \begin{bmatrix} X_{i1} \\ \vdots \\ X_{ip} \end{bmatrix}$$

and ($j = 1, \dots, p$)

$$\mathbf{x}_j = \begin{bmatrix} X_{1j} \\ \vdots \\ X_{nj} \end{bmatrix},$$

we can write

$$\mathbf{X}' = [\mathbf{x}_1 \cdots \mathbf{x}_n]$$

and

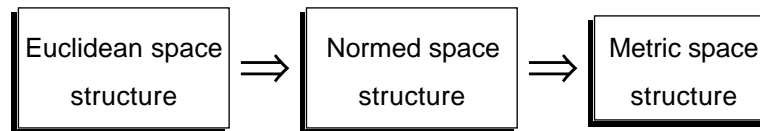
$$\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_p].$$

Considering the notation just introduced, we say that $\mathbf{x}_1, \dots, \mathbf{x}_n$ and $\mathbf{x}_1, \dots, \mathbf{x}_p$ represent, respectively, the n individuals and the p variables.

2.2 ALGEBRAIC STRUCTURE

Regarding $\mathbf{x}_1, \dots, \mathbf{x}_n$ and $\mathbf{x}_1, \dots, \mathbf{x}_p$ as elements of \mathbb{R}^p and \mathbb{R}^n , respectively, \mathbb{R}^p (*individual space*) and \mathbb{R}^n (*variable space*) are equipped with a Euclidean metric.

Obviously, the introduction of a Euclidean metric allows us to calculate in \mathbb{R}^p and \mathbb{R}^n , in addition to the inner product between vectors, both the length of vectors and the distance between vectors.



2.2.1 EUCLIDEAN METRIC IN THE INDIVIDUAL SPACE

In \mathbb{R}^p the matrix (symmetric and positive definite (p.d.)) of the Euclidean metric – with respect to the basis consisting of the p canonical vectors $\mathbf{u}_1, \dots, \mathbf{u}_p$ – is generally of the form

$$\mathbf{Q} = \text{diag}(q_1, \dots, q_p)$$

where $q_j > 0$ ($j = 1, \dots, p$) represents the *weight* given to the j th variable and denotes its «importance» in the set of the p variables ⁽²⁾.

The choice of the weights q_1, \dots, q_p generally depends on the measurement units and/or the variances of the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$.

The situations which may occur are:

- the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ are expressed in the same measurement unit and present approximatively the same variance;
- the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ are expressed in the same measurement unit but present considerably different variances;
- the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ are expressed in different measurement units.

(2) In contexts differing from PCA (e.g., canonical correlation analysis), the metric is specified in other ways.

In the first case, the weights q_1, \dots, q_p are usually chosen setting $q_1 = \dots = q_p = 1$ – namely, assuming that each variable has the same importance of all the others – and thus $\mathbf{Q} = \mathbf{I}_p$.

In the remaining two cases, we often choose the weights q_1, \dots, q_p as the reciprocals of the variances of the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$. The meaning of this choice will be explained below (Section 5.3).

2.2.2 EUCLIDEAN METRIC IN THE VARIABLE SPACE

In \mathbb{R}^n the matrix (symmetric and p.d.) of the Euclidean metric – with respect to the basis consisting of the n canonical vectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ – is

$$\mathbf{M} = \text{diag}(m_1, \dots, m_n)$$

where $m_i > 0$ ($i = 1, \dots, n$), $\sum_i m_i = 1$, represents the *weight* given to the i th individual and denotes its «importance» in the set of the n individuals.

Whenever we do not have sufficient indications about the differing importance of the n individuals, we can set $m_1 = \dots = m_n = m^*$ from which – taking into account the condition $\sum_i m_i = 1$ – we obtain $m^* = \frac{1}{n}$ and thus

$$\mathbf{M} = \text{diag}\left(\frac{1}{n}, \dots, \frac{1}{n}\right).$$

2.3 CENTRED DATA MATRICES

2.3.1 GENERAL CENTRED DATA MATRIX

Given any vector

$$\mathbf{c}_* = \begin{bmatrix} c_{*1} \\ \vdots \\ c_{*p} \end{bmatrix},$$

consider the matrix (*general centred data matrix*)

$$\mathbf{Z} = \mathbf{X} - \mathbf{u} \mathbf{c}_*' = \begin{bmatrix} X_{11} & \cdots & X_{1p} \\ \cdots & \cdots & \cdots \\ X_{n1} & \cdots & X_{np} \end{bmatrix} - \begin{bmatrix} c_{*1} & \cdots & c_{*p} \\ \cdots & \cdots & \cdots \\ c_{*1} & \cdots & c_{*p} \end{bmatrix} = \begin{bmatrix} X_{11} - c_{*1} & \cdots & X_{1p} - c_{*p} \\ \cdots & \cdots & \cdots \\ X_{n1} - c_{*1} & \cdots & X_{np} - c_{*p} \end{bmatrix}$$

where \mathbf{u} is a column vector of order n with elements all equal to 1.

Then, setting ($i = 1, \dots, n$)

$$\mathbf{z}_i = \begin{bmatrix} X_{i1} - c_{*1} \\ \vdots \\ X_{ip} - c_{*p} \end{bmatrix} = \begin{bmatrix} X_{i1} \\ \vdots \\ X_{ip} \end{bmatrix} - \begin{bmatrix} c_{*1} \\ \vdots \\ c_{*p} \end{bmatrix} = \mathbf{x}_i - \mathbf{c}_*$$

and ($j = 1, \dots, p$)

$$\mathbf{z}_j = \begin{bmatrix} X_{1j} - \bar{c}_{*j} \\ \vdots \\ X_{nj} - \bar{c}_{*j} \end{bmatrix} = \begin{bmatrix} X_{1j} \\ \vdots \\ X_{nj} \end{bmatrix} - \begin{bmatrix} \bar{c}_{*j} \\ \vdots \\ \bar{c}_{*j} \end{bmatrix} = \mathbf{x}_j - \bar{\mathbf{c}}_{*j},$$

we can write

$$\mathbf{Z}' = [\mathbf{z}_1 \cdots \mathbf{z}_n]$$

and

$$\mathbf{Z} = [\mathbf{z}_1 \cdots \mathbf{z}_p].$$

2.3.2 MEAN CENTRED DATA MATRIX

Let

$$\mathbf{g} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

where $\bar{x}_j = \sum_i m_i x_{ij}$ is the (weighted) arithmetic mean of the variable \mathbf{x}_j .

Notice that we can write

$$\begin{aligned} \mathbf{g} &= \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} = \begin{bmatrix} \sum_i m_i x_{i1} \\ \vdots \\ \sum_i m_i x_{ip} \end{bmatrix} = \begin{bmatrix} x_{11} \cdots x_{n1} \\ \cdots \cdots \cdots \\ x_{1p} \cdots x_{np} \end{bmatrix} \begin{bmatrix} m_1 \\ \vdots \\ m_n \end{bmatrix} \\ &= [\mathbf{x}_1 \cdots \mathbf{x}_n] \begin{bmatrix} m_1 \\ \vdots \\ m_n \end{bmatrix} = \begin{bmatrix} \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_p \end{bmatrix} \begin{bmatrix} m_1 \\ \vdots \\ m_n \end{bmatrix} = \mathbf{X}'\mathbf{M}\mathbf{u}. \end{aligned}$$

The vector \mathbf{g} is called the *barycentre* (*centroid*) of the n individuals

$\mathbf{x}_1, \dots, \mathbf{x}_n$ or the *mean vector* of the p variables $\mathbf{x}_1, \dots, \mathbf{x}_p$.

Next, consider the matrix (*mean centred data matrix*)

$$\mathbf{Y} = \mathbf{X} - \mathbf{u} \mathbf{g}' = \begin{bmatrix} X_{11} & \cdots & X_{1p} \\ \cdots & \cdots & \cdots \\ X_{n1} & \cdots & X_{np} \end{bmatrix} - \begin{bmatrix} \bar{X}_1 & \cdots & \bar{X}_p \\ \cdots & \cdots & \cdots \\ \bar{X}_1 & \cdots & \bar{X}_p \end{bmatrix} = \begin{bmatrix} X_{11} - \bar{X}_1 & \cdots & X_{1p} - \bar{X}_p \\ \cdots & \cdots & \cdots \\ X_{n1} - \bar{X}_1 & \cdots & X_{np} - \bar{X}_p \end{bmatrix}.$$

Then, setting ($i = 1, \dots, n$)

$$\mathbf{y}_i = \begin{bmatrix} X_{i1} - \bar{X}_1 \\ \vdots \\ X_{ip} - \bar{X}_p \end{bmatrix} = \begin{bmatrix} X_{i1} \\ \vdots \\ X_{ip} \end{bmatrix} - \begin{bmatrix} \bar{X}_1 \\ \vdots \\ \bar{X}_p \end{bmatrix} = \mathbf{x}_i - \mathbf{g}$$

and ($j = 1, \dots, p$)

$$\mathbf{y}_j = \begin{bmatrix} X_{1j} - \bar{X}_j \\ \vdots \\ X_{nj} - \bar{X}_j \end{bmatrix} = \begin{bmatrix} X_{1j} \\ \vdots \\ X_{nj} \end{bmatrix} - \begin{bmatrix} \bar{X}_j \\ \vdots \\ \bar{X}_j \end{bmatrix} = \mathbf{x}_j - \bar{\mathbf{x}}_j,$$

we can write

$$\mathbf{Y}' = [\mathbf{y}_1 \cdots \mathbf{y}_n]$$

and

$$\bar{\mathbf{Y}} = [\mathbf{y}_1 \cdots \mathbf{y}_p].$$

Taking into account the notation just introduced, we say that $\mathbf{y}_1, \dots, \mathbf{y}_n$ and $\mathbf{y}_1, \dots, \mathbf{y}_p$ represent, respectively, the n individuals and the p variables (measured in terms of deviations from the means).

Of course, the (weighted) arithmetic mean of each \mathbf{y}_j ($j = 1, \dots, p$) is zero.

REMARK 1. Notice that

$$\begin{aligned} \mathbf{Z} &= \mathbf{X} - \mathbf{u} \mathbf{c}'_* && = \mathbf{X} - \mathbf{u} (\mathbf{c}'_* + \mathbf{g} - \mathbf{g})' \\ &= \mathbf{X} - \mathbf{u} \mathbf{g}' - \mathbf{u} (\mathbf{c}'_* - \mathbf{g})' && = \mathbf{Y} - \mathbf{u} (\mathbf{c}'_* - \mathbf{g})'. \end{aligned}$$

3 PRELIMINARY CONCEPTS IN THE INDIVIDUAL SPACE

3.1 INERTIA RELATIVE TO A VECTOR

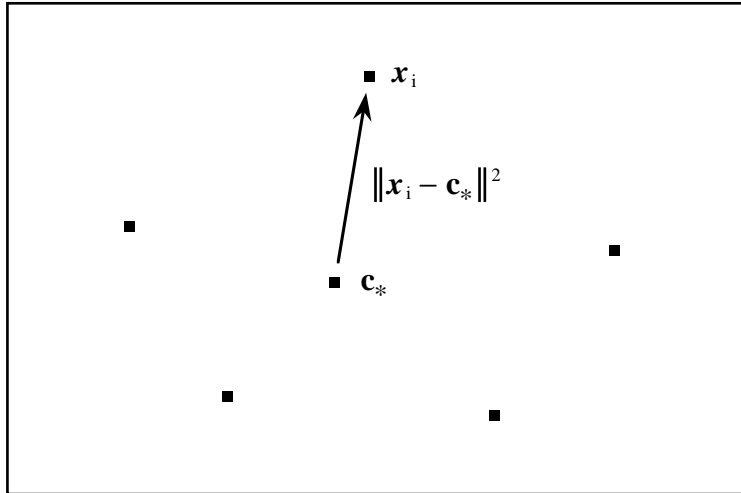
Consider the n individuals $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ with weights given, respectively, by m_1, \dots, m_n and a generic vector $\mathbf{c}_* \in \mathbb{R}^p$.

The quantity

$$I_{\mathbf{c}_*} = \sum_i m_i \|\mathbf{x}_i - \mathbf{c}_*\|^2 = \sum_i m_i (\mathbf{x}_i - \mathbf{c}_*)' \mathbf{Q} (\mathbf{x}_i - \mathbf{c}_*)$$

is called the *inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{c}_** and represents a (weighted) dispersion measure of $\mathbf{x}_1, \dots, \mathbf{x}_n$ with respect to \mathbf{c}_* (Fig. 1)⁽³⁾.

Fig. 1



In turn, the quantity

$$I_{\mathbf{g}} = \sum_i m_i \|\mathbf{x}_i - \mathbf{g}\|^2 = \sum_i m_i (\mathbf{x}_i - \mathbf{g})' \mathbf{Q} (\mathbf{x}_i - \mathbf{g})$$

is called the *inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to the barycentre*.

Notice that, taking into account the notations introduced above (Sections 2.3.1 and 2.3.2), we can write

(3) Some of the concepts introduced in this Section 3 from the point of view of the individual space will be reinterpreted in the following Section 4 from the point of view of the variable space.

$$\begin{aligned}
I_{\mathbf{c}_*} &= \sum_i m_i \mathbf{z}_i' \mathbf{Q} \mathbf{z}_i &&= \text{tr} \left\{ \begin{bmatrix} m_1 \mathbf{z}_1' \mathbf{Q} \mathbf{z}_1 & \cdots & m_1 \mathbf{z}_1' \mathbf{Q} \mathbf{z}_n \\ \cdots & \cdots & \cdots \\ m_n \mathbf{z}_n' \mathbf{Q} \mathbf{z}_1 & \cdots & m_n \mathbf{z}_n' \mathbf{Q} \mathbf{z}_n \end{bmatrix} \right\} \\
&= \text{tr} \left\{ \mathbf{M} \begin{bmatrix} \mathbf{z}_1' \mathbf{Q} \mathbf{z}_1 & \cdots & \mathbf{z}_1' \mathbf{Q} \mathbf{z}_n \\ \cdots & \cdots & \cdots \\ \mathbf{z}_n' \mathbf{Q} \mathbf{z}_1 & \cdots & \mathbf{z}_n' \mathbf{Q} \mathbf{z}_n \end{bmatrix} \right\} &&= \text{tr} \left\{ \mathbf{M} \begin{bmatrix} \mathbf{z}_1' \\ \vdots \\ \mathbf{z}_n' \end{bmatrix} \mathbf{Q} [\mathbf{z}_1 \cdots \mathbf{z}_n] \right\} \\
&= \text{tr} \{ \mathbf{M} \mathbf{Z} \mathbf{Q} \mathbf{Z}' \} &&= \text{tr} \{ \mathbf{Z}' \mathbf{M} \mathbf{Z} \mathbf{Q} \} \\
&= \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} \}
\end{aligned}$$

and, analogously,

$$\begin{aligned}
I_{\mathbf{g}} &= \sum_i m_i \mathbf{y}_i' \mathbf{Q} \mathbf{y}_i = \text{tr} \{ \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \} \\
&= \text{tr} \{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} \}.
\end{aligned}$$

The matrices

$$\mathbf{V}_{\mathbf{c}_*} = [\mathbf{z}_1 \cdots \mathbf{z}_n] \mathbf{M} \begin{bmatrix} \mathbf{z}_1' \\ \vdots \\ \mathbf{z}_n' \end{bmatrix} = \mathbf{Z}' \mathbf{M} \mathbf{Z} \quad , \quad \mathbf{V}_{\mathbf{g}} = [\mathbf{y}_1 \cdots \mathbf{y}_n] \mathbf{M} \begin{bmatrix} \mathbf{y}_1' \\ \vdots \\ \mathbf{y}_n' \end{bmatrix} = \mathbf{Y}' \mathbf{M} \mathbf{Y}$$

denote the so-called *inertia matrices of* $\mathbf{x}_1, \dots, \mathbf{x}_n$ *relative, respectively, to* \mathbf{c}_* *and* \mathbf{g} .

3.1.1 HUIGHENS THEOREM

Taking into account that $\mathbf{Z} = \mathbf{Y} - \mathbf{u}(\mathbf{c}_* - \mathbf{g})'$ (Remark 1) and that, as can easily be verified, $\mathbf{Y}' \mathbf{M} \mathbf{u} = \mathbf{0}$, we have

$$\begin{aligned}
\mathbf{V}_{\mathbf{c}_*} &= \mathbf{Z}' \mathbf{M} \mathbf{Z} \\
&= (\mathbf{Y} - \mathbf{u}(\mathbf{c}_* - \mathbf{g})')' \mathbf{M} (\mathbf{Y} - \mathbf{u}(\mathbf{c}_* - \mathbf{g})') \\
&= \mathbf{Y}' \mathbf{M} \mathbf{Y} - \mathbf{Y}' \mathbf{M} \mathbf{u} (\mathbf{c}_* - \mathbf{g})' - (\mathbf{c}_* - \mathbf{g}) \mathbf{u}' \mathbf{M} \mathbf{Y} + (\mathbf{c}_* - \mathbf{g}) \mathbf{u}' \mathbf{M} \mathbf{u} (\mathbf{c}_* - \mathbf{g})' \\
&= \mathbf{V}_{\mathbf{g}} + (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})'.
\end{aligned}$$

Thus, we get (*Huighens theorem*)

$$\begin{aligned} I_{\mathbf{c}_*} &= \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} \} &&= \text{tr} \{ [\mathbf{V}_{\mathbf{g}} + (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})'] \mathbf{Q} \} \\ &= \text{tr} \{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} \} + (\mathbf{c}_* - \mathbf{g})' \mathbf{Q} (\mathbf{c}_* - \mathbf{g}) = I_{\mathbf{g}} + \| \mathbf{c}_* - \mathbf{g} \|^2 . \end{aligned}$$

In other words, the inertia $I_{\mathbf{c}_*}$ of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{c}_* may be split up into the sum of two addenda:

- $I_{\mathbf{g}}$ which represents the inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{g} ;
- $\| \mathbf{c}_* - \mathbf{g} \|^2$ which represents the square distance between \mathbf{c}_* and \mathbf{g} .

$I_{\mathbf{c}_*}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{c}_*	=	$I_{\mathbf{g}}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{g}	+	$\ \mathbf{c}_* - \mathbf{g} \ ^2$ square distance between \mathbf{c}_* and \mathbf{g}
--	---	--	---	---

REMARK 2. Notice that $I_{\mathbf{c}_*}$ reaches the minimum $I_{\mathbf{g}}$ when $\mathbf{c}_* = \mathbf{g}$.

3.2 INERTIA ALONG A LINEAR VARIETY

Consider in \mathbb{R}^p a subspace C_k of dimension k ($1 \leq k < p$), and its orthogonal complement C_k^\perp .

Denote the orthogonal projection matrices on C_k and C_k^\perp , respectively, by \mathbf{P} and $\mathbf{I}_p - \mathbf{P}$.

Of course, both these matrices are idempotent and selfadjoint, namely

$$\begin{aligned} \mathbf{P}^2 &= \mathbf{P} && , && \mathbf{P}'\mathbf{Q} = \mathbf{Q}\mathbf{P} \\ (\mathbf{I}_p - \mathbf{P})^2 &= (\mathbf{I}_p - \mathbf{P}) && , && (\mathbf{I}_p - \mathbf{P})'\mathbf{Q} = \mathbf{Q}(\mathbf{I}_p - \mathbf{P}) . \end{aligned}$$

Successively, consider the linear variety $\mathbf{c}_* + C_k$ of direction C_k and translation \mathbf{c}_* .

Clearly, the vector $\widehat{\mathbf{x}}_i + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)$ – where $\widehat{\mathbf{x}}_i = \mathbf{P}\mathbf{x}_i$ and $\widehat{\mathbf{c}}_* = \mathbf{P}\mathbf{c}_*$ denote the orthogonal projections, respectively, of \mathbf{x}_i and \mathbf{c}_* on C_k – is the orthogonal projection of \mathbf{x}_i on $\mathbf{c}_* + C_k$ (Fig. 2).

The quantity

$$I_{\mathbf{c}_*+C_k} = \sum_i m_i \left\| (\widehat{\mathbf{x}}_i + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)) - \mathbf{c}_* \right\|^2 = \sum_i m_i \left\| \widehat{\mathbf{x}}_i - \widehat{\mathbf{c}}_* \right\|^2$$

is called the *inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ along \mathbf{c}_*+C_k* or *explained by \mathbf{c}_*+C_k* and represents a (weighted) dispersion measure of the projected vectors $\widehat{\mathbf{x}}_1 + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)$, \dots , $\widehat{\mathbf{x}}_n + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)$ with respect to \mathbf{c}_* .

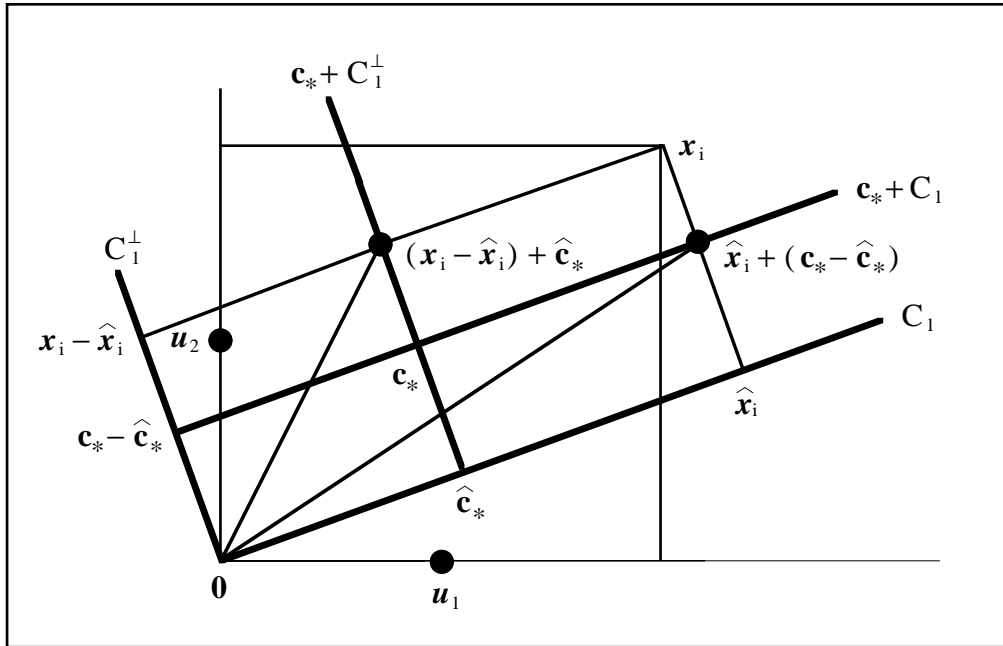
Analogously, the vector $(\mathbf{x}_i - \widehat{\mathbf{x}}_i) + \widehat{\mathbf{c}}_*$ is the orthogonal projection of \mathbf{x}_i on the linear variety $\mathbf{c}_*+C_k^\perp$ of direction C_k^\perp and translation \mathbf{c}_* .

The quantity

$$I_{\mathbf{c}_*+C_k^\perp} = \sum_i m_i \left\| ((\mathbf{x}_i - \widehat{\mathbf{x}}_i) + \widehat{\mathbf{c}}_*) - \mathbf{c}_* \right\|^2 = \sum_i m_i \left\| \mathbf{x}_i - (\widehat{\mathbf{x}}_i + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)) \right\|^2$$

is called the *inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ along $\mathbf{c}_*+C_k^\perp$* or *not explained by \mathbf{c}_*+C_k* and represents a (weighted) dispersion measure of the projected vectors $(\mathbf{x}_1 - \widehat{\mathbf{x}}_1) + \widehat{\mathbf{c}}_*$, \dots , $(\mathbf{x}_n - \widehat{\mathbf{x}}_n) + \widehat{\mathbf{c}}_*$ with respect to \mathbf{c}_* .

Fig. 2



Now, consider the linear variety $\mathbf{g}+C_k$ of direction C_k and translation \mathbf{g} .

The vector $\widehat{\mathbf{x}}_i + (\mathbf{g} - \widehat{\mathbf{g}})$ – where $\widehat{\mathbf{x}}_i = \mathbf{P}\mathbf{x}_i$ and $\widehat{\mathbf{g}} = \mathbf{P}\mathbf{g}$ denote the orthogonal projections, respectively, of \mathbf{x}_i and \mathbf{g} on C_k – is the orthogonal

projection of \mathbf{x}_i on $\mathbf{g} + C_k$ (Fig. 3).

The quantity

$$I_{\mathbf{g}+C_k} = \sum_i m_i \left\| (\hat{\mathbf{x}}_i + (\mathbf{g} - \hat{\mathbf{g}})) - \mathbf{g} \right\|^2 = \sum_i m_i \left\| \hat{\mathbf{x}}_i - \hat{\mathbf{g}} \right\|^2$$

is called the *inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ along $\mathbf{g} + C_k$* or *explained by $\mathbf{g} + C_k$* and represents a (weighted) dispersion measure of the projected vectors $\hat{\mathbf{x}}_1 + (\mathbf{g} - \hat{\mathbf{g}}), \dots, \hat{\mathbf{x}}_n + (\mathbf{g} - \hat{\mathbf{g}})$ with respect to \mathbf{g} .

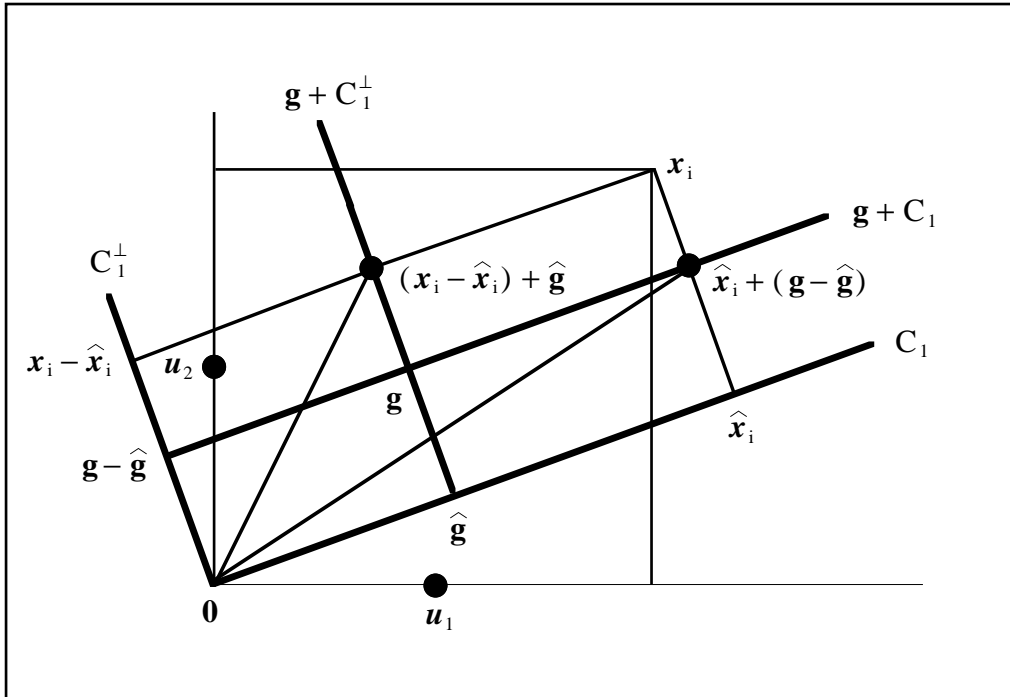
Analogously, the vector $(\mathbf{x}_i - \hat{\mathbf{x}}_i) + \hat{\mathbf{g}}$ is the orthogonal projection of \mathbf{x}_i on the linear variety $\mathbf{g} + C_k^\perp$ of direction C_k^\perp and translation \mathbf{g} .

The quantity

$$I_{\mathbf{g}+C_k^\perp} = \sum_i m_i \left\| ((\mathbf{x}_i - \hat{\mathbf{x}}_i) + \hat{\mathbf{g}}) - \mathbf{g} \right\|^2 = \sum_i m_i \left\| \mathbf{x}_i - (\hat{\mathbf{x}}_i + (\mathbf{g} - \hat{\mathbf{g}})) \right\|^2$$

is called the *inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ along $\mathbf{g} + C_k^\perp$* or *not explained by $\mathbf{g} + C_k$* and represents a (weighted) dispersion measure of the projected vectors $(\mathbf{x}_1 - \hat{\mathbf{x}}_1) + \hat{\mathbf{g}}, \dots, (\mathbf{x}_n - \hat{\mathbf{x}}_n) + \hat{\mathbf{g}}$ with respect to \mathbf{g} .

Fig. 3



3.2.1 A DECOMPOSITION OF THE INERTIA RELATIVE TO A VECTOR

Firstly, notice that we can write

$$\begin{aligned} I_{\mathbf{c}_*} &= \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} \} &= \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} \mathbf{P} + \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} (\mathbf{I}_p - \mathbf{P}) \} \\ &= \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} \mathbf{P} \} + \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} (\mathbf{I}_p - \mathbf{P}) \} . \end{aligned}$$

Thus – since we have $(\mathbf{P}' \mathbf{Q} \mathbf{P} = \mathbf{Q} \mathbf{P})$

$$\begin{aligned} I_{\mathbf{c}_* + \mathbf{C}_k} &= \sum_i m_i \| \mathbf{P} \mathbf{x}_i - \mathbf{P} \mathbf{c}_* \|^2 &= \sum_i m_i \| \mathbf{P} (\mathbf{x}_i - \mathbf{c}_*) \|^2 \\ &= \sum_i m_i \| \mathbf{P} \mathbf{z}_i \|^2 &= \sum_i m_i \mathbf{z}_i' \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{z}_i \\ &= \text{tr} \left\{ \begin{bmatrix} m_1 \mathbf{z}_1' \mathbf{Q} \mathbf{P} \mathbf{z}_1 & \cdots & m_1 \mathbf{z}_1' \mathbf{Q} \mathbf{P} \mathbf{z}_n \\ \cdots & \cdots & \cdots \\ m_n \mathbf{z}_n' \mathbf{Q} \mathbf{P} \mathbf{z}_1 & \cdots & m_n \mathbf{z}_n' \mathbf{Q} \mathbf{P} \mathbf{z}_n \end{bmatrix} \right\} &= \text{tr} \left\{ \mathbf{M} \begin{bmatrix} \mathbf{z}_1' \\ \vdots \\ \mathbf{z}_n' \end{bmatrix} \mathbf{Q} \mathbf{P} \begin{bmatrix} \mathbf{z}_1 \\ \cdots \\ \mathbf{z}_n \end{bmatrix} \right\} \\ &= \text{tr} \{ \mathbf{M} \mathbf{Z} \mathbf{Q} \mathbf{P} \mathbf{Z}' \} &= \text{tr} \{ \mathbf{Z}' \mathbf{M} \mathbf{Z} \mathbf{Q} \mathbf{P} \} \\ &= \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} \mathbf{P} \} \end{aligned}$$

and, analogously,

$$I_{\mathbf{c}_* + \mathbf{C}_k^\perp} = \text{tr} \{ \mathbf{V}_{\mathbf{c}_*} \mathbf{Q} (\mathbf{I}_p - \mathbf{P}) \}$$

– we obtain the decomposition

$$I_{\mathbf{c}_*} = I_{\mathbf{c}_* + \mathbf{C}_k} + I_{\mathbf{c}_* + \mathbf{C}_k^\perp}$$

Namely, the inertia $I_{\mathbf{c}_*}$ of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{c}_* may be split up into the sum of two addenda:

- $I_{\mathbf{c}_* + \mathbf{C}_k}$ which is the inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ explained by $\mathbf{c}_* + \mathbf{C}_k$;
- $I_{\mathbf{c}_* + \mathbf{C}_k^\perp}$ which is the inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ not explained by $\mathbf{c}_* + \mathbf{C}_k$.

$I_{\mathbf{c}_*}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ relative to \mathbf{c}_*	=	$I_{\mathbf{c}_* + \mathbf{C}_k}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ explained by $\mathbf{c}_* + \mathbf{C}_k$	+	$I_{\mathbf{c}_* + \mathbf{C}_k^\perp}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ not explained by $\mathbf{c}_* + \mathbf{C}_k$
--	---	---	---	---

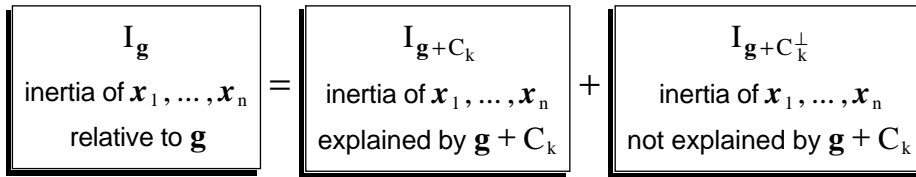
Of course, for $\mathbf{c}_* = \mathbf{g}$, we get

$$I_{\mathbf{g}} = I_{\mathbf{g}+C_k} + I_{\mathbf{g}+C_k^\perp}$$

where $(\hat{\mathbf{g}} = \mathbf{P}\mathbf{g})$

$$I_{\mathbf{g}+C_k} = \sum_i m_i \|\hat{\mathbf{x}}_i - \hat{\mathbf{g}}\|^2 = \text{tr} \{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} \mathbf{P} \} ,$$

$$I_{\mathbf{g}+C_k^\perp} = \sum_i m_i \|(x_i - \hat{x}_i) + \hat{\mathbf{g}} - \mathbf{g}\|^2 = \text{tr} \{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} (\mathbf{I}_p - \mathbf{P}) \} .$$

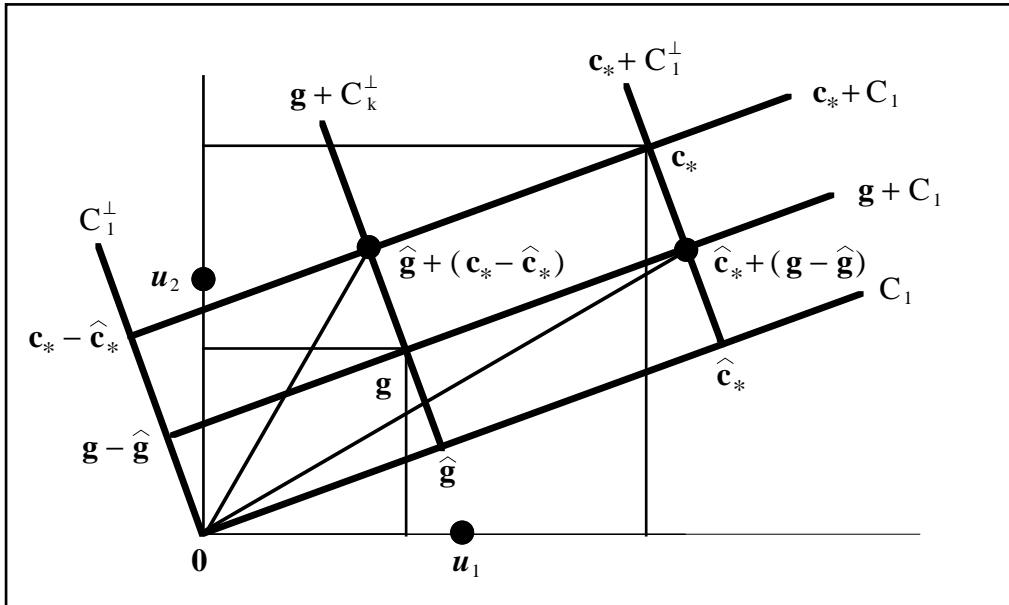


3.2.2 A DECOMPOSITION OF THE INERTIA ALONG A LINEAR VARIETY

Firstly, consider the linear varieties \mathbf{c}_*+C_k and $\mathbf{g}+C_k$. The vector $\hat{\mathbf{c}}_*+(\mathbf{g}-\hat{\mathbf{g}})$ is the orthogonal projection of \mathbf{c}_* on the linear variety $\mathbf{g}+C_k$ (Fig. 4). Hence, the square distance between \mathbf{c}_*+C_k and $\mathbf{g}+C_k$ is

$$\|\mathbf{c}_* - \hat{\mathbf{c}}_* - (\mathbf{g} - \hat{\mathbf{g}})\|^2 = \|(\mathbf{c}_* - \mathbf{g}) - (\hat{\mathbf{c}}_* - \hat{\mathbf{g}})\|^2 .$$

Fig. 4



Successively, consider the linear varieties $\mathbf{c}_* + C_k^\perp$ and $\mathbf{g} + C_k^\perp$. The vector $\widehat{\mathbf{g}} + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)$ is the orthogonal projection of \mathbf{c}_* on the linear variety $\mathbf{g} + C_k^\perp$. Therefore, the square distance between $\mathbf{c}_* + C_k^\perp$ and $\mathbf{g} + C_k^\perp$ is

$$\|\mathbf{c}_* - \widehat{\mathbf{g}} - (\mathbf{c}_* - \widehat{\mathbf{c}}_*)\|^2 = \|\widehat{\mathbf{c}}_* - \widehat{\mathbf{g}}\|^2.$$

Now, notice that – since

$$I_{\mathbf{c}_* + C_k} = I_{\mathbf{c}_*} - I_{\mathbf{c}_* + C_k^\perp}, \quad I_{\mathbf{g} + C_k} = I_{\mathbf{g}} - I_{\mathbf{g} + C_k^\perp}, \quad \mathbf{V}_{\mathbf{c}_*} = \mathbf{V}_{\mathbf{g}} + (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})'$$

– we get ($\mathbf{P}'\mathbf{Q}\mathbf{P} = \mathbf{Q}\mathbf{P}$)

$$\begin{aligned} I_{\mathbf{c}_* + C_k} &= \text{tr} \left\{ [\mathbf{V}_{\mathbf{g}} + (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})'] \mathbf{Q} \right\} - \text{tr} \left\{ [\mathbf{V}_{\mathbf{g}} + (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})'] \mathbf{Q} (\mathbf{I} - \mathbf{P}) \right\} \\ &= \text{tr} \left\{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} \right\} - \text{tr} \left\{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} (\mathbf{I}_p - \mathbf{P}) \right\} + \text{tr} \left\{ (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})' \mathbf{Q} \right\} \\ &\quad - \text{tr} \left\{ (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})' \mathbf{Q} (\mathbf{I}_p - \mathbf{P}) \right\} \\ &= \text{tr} \left\{ \mathbf{V}_{\mathbf{g}} \mathbf{Q} \mathbf{P} \right\} + \text{tr} \left\{ (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})' \mathbf{Q} \mathbf{P} \right\} \\ &= I_{\mathbf{g} + C_k} + \text{tr} \left\{ (\mathbf{c}_* - \mathbf{g})(\mathbf{c}_* - \mathbf{g})' \mathbf{P}' \mathbf{Q} \mathbf{P} \right\} \\ &= I_{\mathbf{g} + C_k} + (\mathbf{c}_* - \mathbf{g})' \mathbf{P}' \mathbf{Q} \mathbf{P} (\mathbf{c}_* - \mathbf{g}) \\ &= I_{\mathbf{g} + C_k} + (\mathbf{P} \mathbf{c}_* - \mathbf{P} \mathbf{g})' \mathbf{Q} (\mathbf{P} \mathbf{c}_* - \mathbf{P} \mathbf{g}) \\ &= I_{\mathbf{g} + C_k} + (\widehat{\mathbf{c}}_* - \widehat{\mathbf{g}})' \mathbf{Q} (\widehat{\mathbf{c}}_* - \widehat{\mathbf{g}}) \\ &= I_{\mathbf{g} + C_k} + \|\widehat{\mathbf{c}}_* - \widehat{\mathbf{g}}\|^2. \end{aligned}$$

Namely, the inertia $I_{\mathbf{c}_* + C_k}$ of $\mathbf{x}_1, \dots, \mathbf{x}_n$ explained by $\mathbf{c}_* + C_k$ may be split up into the sum of two addenda:

- $I_{\mathbf{g} + C_k}$ which is the inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ explained by $\mathbf{g} + C_k$;
- the square distance between the linear varieties $\mathbf{c}_* + C_k^\perp$ and $\mathbf{g} + C_k^\perp$.

$$\boxed{\begin{array}{c} I_{\mathbf{c}_* + C_k} \\ \text{inertia of } \mathbf{x}_1, \dots, \mathbf{x}_n \\ \text{explained by } \mathbf{c}_* + C_k \end{array}} = \boxed{\begin{array}{c} I_{\mathbf{g} + C_k} \\ \text{inertia of } \mathbf{x}_1, \dots, \mathbf{x}_n \\ \text{explained by } \mathbf{g} + C_k \end{array}} + \boxed{\begin{array}{c} \text{square distance} \\ \text{between} \\ \mathbf{c}_* + C_k^\perp, \mathbf{g} + C_k^\perp \end{array}}$$

Analogously, as can easily be verified, we obtain the decomposition

$$I_{\mathbf{c}_*+C_k^\perp} = I_{\mathbf{g}+C_k^\perp} + \|\mathbf{c}_* - \widehat{\mathbf{c}}_* - (\mathbf{g} - \widehat{\mathbf{g}})\|^2$$

which shows that the inertia $I_{\mathbf{c}_*+C_k^\perp}$ of $\mathbf{x}_1, \dots, \mathbf{x}_n$ not explained by \mathbf{c}_*+C_k may be split up into the sum of two addenda:

- $I_{\mathbf{g}+C_k^\perp}$ which is the inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ not explained by $\mathbf{g}+C_k$;
- the square distance between the linear varieties \mathbf{c}_*+C_k and $\mathbf{g}+C_k$.

$I_{\mathbf{c}_*+C_k^\perp}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ not explained by \mathbf{c}_*+C_k	=	$I_{\mathbf{g}+C_k^\perp}$ inertia of $\mathbf{x}_1, \dots, \mathbf{x}_n$ not explained by $\mathbf{g}+C_k$	+	square distance between $\mathbf{c}_*+C_k, \mathbf{g}+C_k$
---	---	---	---	--

REMARK 3. Notice that both the inertias $I_{\mathbf{c}_*+C_k}$ and $I_{\mathbf{c}_*+C_k^\perp}$ of $\mathbf{x}_1, \dots, \mathbf{x}_n$ are minimized whenever the linear varieties \mathbf{c}_*+C_k and $\mathbf{c}_*+C_k^\perp$ pass through the barycentre \mathbf{g} .

REMARK 4. It is immediately apparent that, if we consider the n individuals $\mathbf{y}_1, \dots, \mathbf{y}_n$ (measured in terms of deviations from the means) instead of $\mathbf{x}_1, \dots, \mathbf{x}_n$, we can interpret:

- $I_{\mathbf{g}}$ as the inertia of $\mathbf{y}_1, \dots, \mathbf{y}_n$ relative to $\mathbf{0}$;
- $I_{\mathbf{g}+C_k}$ as the inertia of $\mathbf{y}_1, \dots, \mathbf{y}_n$ explained by the subspace C_k ;
- $I_{\mathbf{g}+C_k^\perp}$ as the inertia of $\mathbf{y}_1, \dots, \mathbf{y}_n$ not explained by the subspace C_k .

4 PRELIMINARY CONCEPTS IN THE VARIABLE SPACE

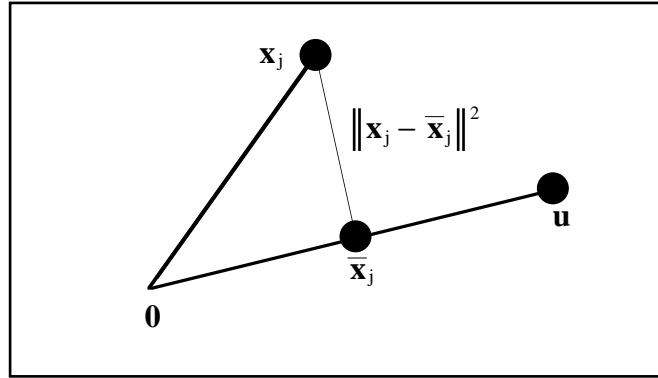
4.1 VARIANCES AND COVARIANCES

Consider the p variables $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$ with weights given, respectively, by q_1, \dots, q_p .

The orthogonal projection $\bar{\mathbf{x}}_j$ of \mathbf{x}_j on the subspace (of dimension 1) spanned by the vector $\mathbf{u} \in \mathbb{R}^n$ with elements all equal to 1 is (Fig. 5)

$$\bar{\mathbf{x}}_j = \mathbf{u} (\mathbf{u}' \mathbf{M} \mathbf{u})^{-1} \mathbf{u}' \mathbf{M} \mathbf{x}_j = \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{x}_j = \mathbf{u} \bar{\mathbf{x}}_j.$$

Fig. 5



The quantity

$$\sigma_j^2 = \|\mathbf{x}_j - \bar{\mathbf{x}}_j\|^2 = (\mathbf{x}_j - \bar{\mathbf{x}}_j)' \mathbf{M} (\mathbf{x}_j - \bar{\mathbf{x}}_j) = \mathbf{y}_j' \mathbf{M} \mathbf{y}_j$$

is the *variance* of \mathbf{x}_j or \mathbf{y}_j and the quantity ($j, t = 1, \dots, p$)

$$\sigma_{jt} = (\mathbf{x}_j - \bar{\mathbf{x}}_j)' \mathbf{M} (\mathbf{x}_t - \bar{\mathbf{x}}_t) = \mathbf{y}_j' \mathbf{M} \mathbf{y}_t$$

is the *covariance* between \mathbf{x}_j and \mathbf{x}_t or between \mathbf{y}_j and \mathbf{y}_t .

In turn, the (symmetric) matrix

$$\mathbf{V} = \begin{bmatrix} \sigma_1^2 & \cdots & \sigma_{1p} \\ \cdots & \cdots & \cdots \\ \sigma_{p1} & \cdots & \sigma_p^2 \end{bmatrix}$$

denotes the so-called *covariance matrix* of $\mathbf{x}_1, \dots, \mathbf{x}_p$ or $\mathbf{y}_1, \dots, \mathbf{y}_p$.

Notice that we can write

$$\mathbf{V} = \begin{bmatrix} \mathbf{y}'_1 \mathbf{M} \mathbf{y}_1 & \cdots & \mathbf{y}'_1 \mathbf{M} \mathbf{y}_p \\ \cdots & \cdots & \cdots \\ \mathbf{y}'_p \mathbf{M} \mathbf{y}_1 & \cdots & \mathbf{y}'_p \mathbf{M} \mathbf{y}_p \end{bmatrix} = \mathbf{Y}' \mathbf{M} \mathbf{Y}.$$

Hence, the covariance matrix \mathbf{V} is nothing other than the inertia matrix \mathbf{V}_g defined above (Section 3.1).

Moreover, we immediately realize that \mathbf{V} is the Gram matrix of $\mathbf{y}_1, \dots, \mathbf{y}_p$ and hence \mathbf{V} is positive definite or positive semi-definite according to $\mathbf{y}_1, \dots, \mathbf{y}_p$ are linearly independent or dependent ⁽⁴⁾.

Of course,

$$r(\mathbf{Y}' \mathbf{M} \mathbf{Y}) = r(\mathbf{Y}' \mathbf{M}^{1/2} \mathbf{M}^{1/2} \mathbf{Y}) = r(\mathbf{M}^{1/2} \mathbf{Y}) = r(\mathbf{Y}).$$

Finally, the quantity

$$J_p = \sum_j q_j \sigma_j^2 = \sum_j q_j \mathbf{y}'_j \mathbf{M} \mathbf{y}_j$$

denotes the so-called *global variability* of $\mathbf{x}_1, \dots, \mathbf{x}_p$ or $\mathbf{y}_1, \dots, \mathbf{y}_p$.

4.2 CORRELATIONS

The cosine of the angle formed by the vectors $\mathbf{y}_j = \mathbf{x}_j - \bar{\mathbf{x}}_j$ and $\mathbf{y}_t = \mathbf{x}_t - \bar{\mathbf{x}}_t$ ($j, t = 1, \dots, p; \mathbf{y}_j, \mathbf{y}_t \neq \mathbf{0}$) is, as can easily be verified, the *linear correlation coefficient* r_{jt} between \mathbf{x}_j and \mathbf{x}_t or between \mathbf{y}_j and \mathbf{y}_t ; namely, we have

$$\cos(\mathbf{y}_j, \mathbf{y}_t) = \frac{(\mathbf{x}_j - \bar{\mathbf{x}}_j)' \mathbf{M} (\mathbf{x}_t - \bar{\mathbf{x}}_t)}{\sigma_j \sigma_t} = \frac{\sigma_{jt}}{\sigma_j \sigma_t} = r_{jt}.$$

In turn, the matrix

$$\mathbf{R} = \mathbf{Q}_{1/\sigma} \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q}_{1/\sigma},$$

where

(4) Notice that, in general, $r(\mathbf{Y}) \neq r(\mathbf{X})$.

$$\mathbf{Q}_{1/\sigma} = \text{diag} \left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_p} \right),$$

represents the so-called *correlation matrix* of $\mathbf{x}_1, \dots, \mathbf{x}_p$ or $\mathbf{y}_1, \dots, \mathbf{y}_p$.

Finally, consider a generic variable $\mathbf{x} \in \mathbb{R}^n$ and the corresponding vector $\mathbf{y} = \mathbf{x} - \mathbf{u}\bar{x}$, where $\bar{x} = \sum_i m_i x_i$ is the (weighted) arithmetic mean of \mathbf{x} .

The orthogonal projection $\hat{\mathbf{y}}$ of \mathbf{y} on the subspace spanned by the p vectors $\mathbf{y}_1, \dots, \mathbf{y}_p$, assuming that $r(\mathbf{Y}) = p$, is given by

$$\hat{\mathbf{y}} = \mathbf{Y}(\mathbf{Y}'\mathbf{M}\mathbf{Y})^{-1}\mathbf{Y}'\mathbf{M}\mathbf{y} = \mathbf{P}_Y\mathbf{y}.$$

The square cosine of the angle formed by the vectors \mathbf{y} and $\hat{\mathbf{y}}$ denotes the *square multiple linear correlation coefficient* (*linear determination coefficient*) ρ between \mathbf{x} and $\mathbf{x}_1, \dots, \mathbf{x}_p$ or between \mathbf{y} and $\mathbf{y}_1, \dots, \mathbf{y}_p$.

In fact, since $(\mathbf{M}\mathbf{P}_Y = \mathbf{P}_Y'\mathbf{M})$

$$\mathbf{y}'\mathbf{M}\hat{\mathbf{y}} = \mathbf{y}'\mathbf{M}\mathbf{P}_Y\mathbf{y} = \mathbf{y}'\mathbf{M}\mathbf{P}_Y\mathbf{P}_Y\mathbf{y} = \mathbf{y}'\mathbf{P}_Y'\mathbf{M}\mathbf{P}_Y\mathbf{y} = \hat{\mathbf{y}}'\mathbf{M}\hat{\mathbf{y}},$$

we can write

$$\cos^2(\mathbf{y}, \hat{\mathbf{y}}) = \frac{(\mathbf{y}'\mathbf{M}\hat{\mathbf{y}})^2}{(\mathbf{y}'\mathbf{M}\mathbf{y})(\hat{\mathbf{y}}'\mathbf{M}\hat{\mathbf{y}})} = \frac{(\hat{\mathbf{y}}'\mathbf{M}\hat{\mathbf{y}})^2}{(\mathbf{y}'\mathbf{M}\mathbf{y})(\hat{\mathbf{y}}'\mathbf{M}\hat{\mathbf{y}})} = \frac{\hat{\mathbf{y}}'\mathbf{M}\hat{\mathbf{y}}}{\mathbf{y}'\mathbf{M}\mathbf{y}} = \rho.$$

4.3 INTERPRETATION OF SOME CONCEPTS OF INERTIA

4.3.1 INERTIA RELATIVE TO THE BARYCENTRE

We want to show that the inertia I_g (Section 3.1) is nothing other than the global variability J_p defined above (Section 4.1).

In fact,

$$\begin{aligned} I_g &= \text{tr} \{ \mathbf{V}\mathbf{Q} \} &&= \text{tr} \{ \mathbf{Q}\mathbf{V} \} \\ &= \text{tr} \left\{ \mathbf{Q} \begin{bmatrix} \mathbf{y}'_1\mathbf{M}\mathbf{y}_1 & \cdots & \mathbf{y}'_1\mathbf{M}\mathbf{y}_p \\ \cdots & \cdots & \cdots \\ \mathbf{y}'_p\mathbf{M}\mathbf{y}_1 & \cdots & \mathbf{y}'_p\mathbf{M}\mathbf{y}_p \end{bmatrix} \right\} &&= \text{tr} \left\{ \begin{bmatrix} q_1 \mathbf{y}'_1\mathbf{M}\mathbf{y}_1 & \cdots & q_1 \mathbf{y}'_1\mathbf{M}\mathbf{y}_p \\ \cdots & \cdots & \cdots \\ q_p \mathbf{y}'_p\mathbf{M}\mathbf{y}_1 & \cdots & q_p \mathbf{y}'_p\mathbf{M}\mathbf{y}_p \end{bmatrix} \right\} \\ &= \sum_j q_j \mathbf{y}'_j\mathbf{M}\mathbf{y}_j &&= J_p. \end{aligned}$$

4.3.2 INERTIA ALONG A LINEAR VARIETY THROUGH THE BARYCENTRE

Consider again the n individuals $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ and the orthogonal projection matrix \mathbf{P} on C_k .

Let

$$\mathbf{P}\mathbf{X}' = \mathbf{P}[\mathbf{x}_1 \cdots \mathbf{x}_n] = [\mathbf{P}\mathbf{x}_1 \cdots \mathbf{P}\mathbf{x}_n] = [\hat{\mathbf{x}}_1 \cdots \hat{\mathbf{x}}_n] = \begin{bmatrix} \hat{x}_{11} & \cdots & \hat{x}_{n1} \\ \vdots & \ddots & \vdots \\ \hat{x}_{1p} & \cdots & \hat{x}_{np} \end{bmatrix}.$$

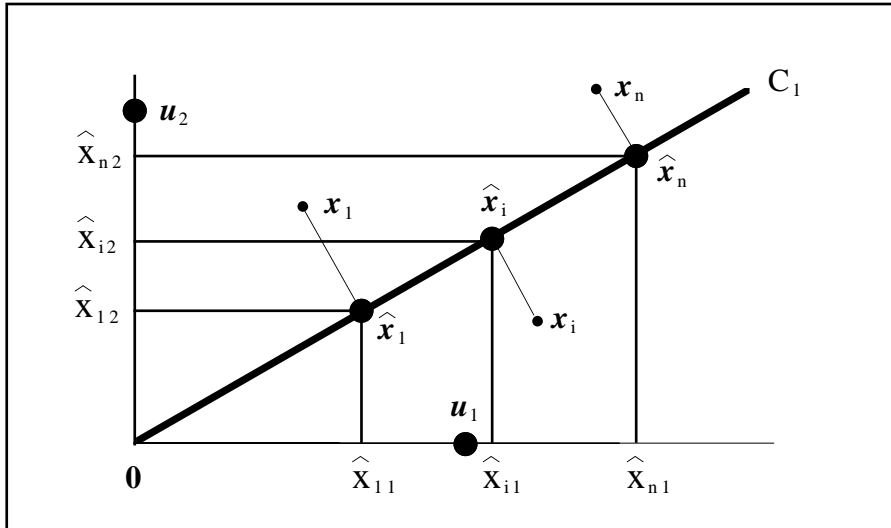
Notice that $\hat{x}_{i1}, \dots, \hat{x}_{ip}$ represent the co-ordinates of $\hat{\mathbf{x}}_i$ (Fig. 6).

Moreover, we can write

$$\mathbf{X}\mathbf{P}' = \begin{bmatrix} \hat{x}_{11} & \cdots & \hat{x}_{1p} \\ \vdots & \ddots & \vdots \\ \hat{x}_{n1} & \cdots & \hat{x}_{np} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_{11} & \cdots & \mathbf{x}_{1p} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_{n1} & \cdots & \mathbf{x}_{np} \end{bmatrix} \begin{bmatrix} \mathbf{p}_{11} & \cdots & \mathbf{p}_{p1} \\ \vdots & \ddots & \vdots \\ \mathbf{p}_{1p} & \cdots & \mathbf{p}_{pp} \end{bmatrix} = [\mathbf{X}\mathbf{p}_1 \cdots \mathbf{X}\mathbf{p}_p]$$

where \mathbf{p}_j ($j = 1, \dots, p$) denotes the j th column of \mathbf{P}' .

Fig. 6



Now, consider the variable ($j = 1, \dots, p$)

$$\hat{\mathbf{x}}_j = \mathbf{X}\mathbf{p}_j = \begin{bmatrix} \hat{x}_{1j} \\ \vdots \\ \hat{x}_{nj} \end{bmatrix}$$

and notice that its (weighted) arithmetic mean and variance are, respectively,

$$\widehat{\mathbf{x}}_j = [m_1 \cdots m_n] \begin{bmatrix} \widehat{x}_{1j} \\ \vdots \\ \widehat{x}_{nj} \end{bmatrix} = \mathbf{u}' \mathbf{M} \widehat{\mathbf{x}}_j = \mathbf{u}' \mathbf{M} \mathbf{X} \mathbf{p}_j = \mathbf{g}' \mathbf{p}_j$$

and

$$\begin{aligned} \widehat{\sigma}_j^2 &= (\mathbf{X} \mathbf{p}_j - \mathbf{u} \mathbf{g}' \mathbf{p}_j)' \mathbf{M} (\mathbf{X} \mathbf{p}_j - \mathbf{u} \mathbf{g}' \mathbf{p}_j) \\ &= \mathbf{p}_j' (\mathbf{X} - \mathbf{u} \mathbf{g}')' \mathbf{M} (\mathbf{X} - \mathbf{u} \mathbf{g}') \mathbf{p}_j \\ &= \mathbf{p}_j' ([\mathbf{x}_1 \cdots \mathbf{x}_p] - \mathbf{u} [\bar{x}_1 \cdots \bar{x}_p])' \mathbf{M} ([\mathbf{x}_1 \cdots \mathbf{x}_p] - \mathbf{u} [\bar{x}_1 \cdots \bar{x}_p]) \mathbf{p}_j \\ &= \mathbf{p}_j' [\mathbf{x}_1 - \mathbf{u} \bar{x}_1 \cdots \mathbf{x}_p - \mathbf{u} \bar{x}_p]' \mathbf{M} [\mathbf{x}_1 - \mathbf{u} \bar{x}_1 \cdots \mathbf{x}_p - \mathbf{u} \bar{x}_p] \mathbf{p}_j \\ &= \mathbf{p}_j' [\mathbf{y}_1 \cdots \mathbf{y}_p]' \mathbf{M} [\mathbf{y}_1 \cdots \mathbf{y}_p] \mathbf{p}_j \\ &= \mathbf{p}_j' \mathbf{V} \mathbf{p}_j. \end{aligned}$$

The quantity

$$\widehat{J}_p = \sum_j q_j \widehat{\sigma}_j^2 = \sum_j q_j \mathbf{p}_j' \mathbf{V} \mathbf{p}_j$$

denotes the *global variability* of $\widehat{\mathbf{x}}_1, \dots, \widehat{\mathbf{x}}_p$.

We want to show that the inertia $I_{\mathbf{g} + C_k}$ of $\mathbf{x}_1, \dots, \mathbf{x}_n$ explained by the linear variety $\mathbf{g} + C_k$ is nothing other than the global variability \widehat{J}_p .

In fact,

$$\begin{aligned} I_{\mathbf{g} + C_k} &= \text{tr} \{ \mathbf{V} \mathbf{Q} \mathbf{P} \} &&= \text{tr} \{ \mathbf{V} \mathbf{P}' \mathbf{Q} \mathbf{P} \} \\ &= \text{tr} \{ \mathbf{Q} \mathbf{P} \mathbf{V} \mathbf{P}' \} &&= \text{tr} \left\{ \mathbf{Q} \begin{bmatrix} \mathbf{p}_1' \\ \vdots \\ \mathbf{p}_p' \end{bmatrix} \mathbf{V} [\mathbf{p}_1 \cdots \mathbf{p}_p] \right\} \\ &= \text{tr} \left\{ \mathbf{Q} \begin{bmatrix} \mathbf{p}_1' \mathbf{V} \mathbf{p}_1 & \cdots & \mathbf{p}_1' \mathbf{V} \mathbf{p}_p \\ \cdots & \cdots & \cdots \\ \mathbf{p}_p' \mathbf{V} \mathbf{p}_1 & \cdots & \mathbf{p}_p' \mathbf{V} \mathbf{p}_p \end{bmatrix} \right\} &&= \text{tr} \left\{ \begin{bmatrix} q_1 \mathbf{p}_1' \mathbf{V} \mathbf{p}_1 & \cdots & q_1 \mathbf{p}_1' \mathbf{V} \mathbf{p}_p \\ \cdots & \cdots & \cdots \\ q_p \mathbf{p}_p' \mathbf{V} \mathbf{p}_1 & \cdots & q_p \mathbf{p}_p' \mathbf{V} \mathbf{p}_p \end{bmatrix} \right\} \\ &= \sum_j q_j \mathbf{p}_j' \mathbf{V} \mathbf{p}_j &&= \widehat{J}_p. \end{aligned}$$

5 AN APPROACH TO PCA

5.1 PRINCIPAL VECTORS, PRINCIPAL COMPONENTS

Consider again the n individuals $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ with weights given, respectively, by m_1, \dots, m_n and the linear varieties $\mathbf{c}_* + C_k$ and $\mathbf{c}_* + C_k^\perp$.

The vectors $\widehat{\mathbf{x}}_1 + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)$, \dots , $\widehat{\mathbf{x}}_n + (\mathbf{c}_* - \widehat{\mathbf{c}}_*)$ may be interpreted as the «images» of $\mathbf{x}_1, \dots, \mathbf{x}_n$ on $\mathbf{c}_* + C_k$.

If we want such images to be, on the whole, the most representative of $\mathbf{x}_1, \dots, \mathbf{x}_n$, a criterion may consist in maximizing the inertia explained by $\mathbf{c}_* + C_k$ with respect to \mathbf{c}_* and C_k ⁽⁵⁾.

This problem can be solved in two steps: at the first step, taking into account Remark 3, we force the linear variety $\mathbf{c}_* + C_k$ to pass through the barycentre \mathbf{g} ; at the second step, we maximize the inertia explained by $\mathbf{g} + C_k$ with respect to C_k .

As regards this last problem, first notice that we may suppose that the orthogonal subspaces C_k and C_k^\perp of \mathbb{R}^p are spanned, respectively, by the orthonormal vectors $\mathbf{c}_1, \dots, \mathbf{c}_k$ and $\mathbf{c}_{k+1}, \dots, \mathbf{c}_p$.

Thus – setting

$$\mathbf{C}_k = [\mathbf{c}_1 \cdots \mathbf{c}_k] \quad , \quad \mathbf{C}_{p-k} = [\mathbf{c}_{k+1} \cdots \mathbf{c}_p] \quad , \quad \mathbf{C}_p = [\mathbf{C}_k \quad \mathbf{C}_{p-k}]$$

– the orthogonal projection matrices on C_k and C_k^\perp become, respectively,

$$\begin{aligned} \mathbf{P} &= \mathbf{C}_k (\mathbf{C}_k' \mathbf{Q} \mathbf{C}_k)^{-1} \mathbf{C}_k' \mathbf{Q} = \mathbf{C}_k \mathbf{C}_k' \mathbf{Q} \\ \mathbf{I}_p - \mathbf{P} &= \mathbf{C}_{p-k} (\mathbf{C}_{p-k}' \mathbf{Q} \mathbf{C}_{p-k})^{-1} \mathbf{C}_{p-k}' \mathbf{Q} = \mathbf{C}_{p-k} \mathbf{C}_{p-k}' \mathbf{Q} . \end{aligned}$$

Moreover,

$$\mathbf{I}_{\mathbf{g} + C_k} = \text{tr} \{ \mathbf{V} \mathbf{Q} \mathbf{P} \} = \text{tr} \{ \mathbf{V} \mathbf{Q} \mathbf{C}_k \mathbf{C}_k' \mathbf{Q} \} = \text{tr} \{ \mathbf{C}_k' \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \} .$$

Hence, our problem lies in finding out

(5) Of course, given that the inertia relative to \mathbf{c}_* is a fixed quantity, this criterion is equivalent to minimizing the inertia not explained by $\mathbf{c}_* + C_k$.

$$\text{Max}_{\mathbf{C}_k} \text{tr} \{ \mathbf{C}'_k \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \} , \quad \mathbf{C}'_k \mathbf{Q} \mathbf{C}_k = \mathbf{I}_k .$$

In order to solve the problem at hand, consider the Lagrange function

$$L(\mathbf{C}_k, \mathbf{L}) = \text{tr} \{ \mathbf{C}'_k \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \} - \text{tr} \{ (\mathbf{C}'_k \mathbf{Q} \mathbf{C}_k - \mathbf{I}_k) \mathbf{L} \}$$

where $\mathbf{L} = \mathbf{L}'$ is a matrix of order (k, k) of Lagrange multipliers.

At $(\tilde{\mathbf{C}}_k, \tilde{\mathbf{L}})$ where $L(\mathbf{C}_k, \mathbf{L})$ has a maximum, it must be

$$\left. \frac{\partial L(\mathbf{C}_k, \mathbf{L})}{\partial \mathbf{C}_k} \right|_{(\tilde{\mathbf{C}}_k, \tilde{\mathbf{L}})} = 2\mathbf{Q}(\mathbf{V}\mathbf{Q}\tilde{\mathbf{C}}_k - \tilde{\mathbf{C}}_k \tilde{\mathbf{L}}) = \mathbf{O}_{(p, k)}$$

$$\left. \frac{\partial L(\mathbf{C}_k, \mathbf{L})}{\partial \mathbf{L}} \right|_{(\tilde{\mathbf{C}}_k, \tilde{\mathbf{L}})} = -(\tilde{\mathbf{C}}_k' \mathbf{Q} \tilde{\mathbf{C}}_k - \mathbf{I}_k) = \mathbf{O}_{(k, k)}$$

which gives

$$\mathbf{V}\mathbf{Q}\tilde{\mathbf{C}}_k = \tilde{\mathbf{C}}_k \tilde{\mathbf{L}} , \quad \tilde{\mathbf{C}}_k' \mathbf{Q} \tilde{\mathbf{C}}_k = \mathbf{I}_k .$$

Therefore, we must look for solutions of the system

$$(*) \quad \mathbf{V}\mathbf{Q}\mathbf{C}_k = \mathbf{C}_k \mathbf{L} , \quad \mathbf{C}'_k \mathbf{Q} \mathbf{C}_k = \mathbf{I}_k$$

in the unknowns \mathbf{C}_k and \mathbf{L} .

To this end, consider the equation

$$\mathbf{V}\mathbf{Q}\mathbf{c} = \lambda \mathbf{c}$$

in the unknowns \mathbf{c} and λ .

This equation possesses p orthonormal eigenvectors $\tilde{\mathbf{c}}_1, \dots, \tilde{\mathbf{c}}_k, \tilde{\mathbf{c}}_{k+1}, \dots, \tilde{\mathbf{c}}_p$ corresponding to the p (real) eigenvalues $\tilde{\lambda}_1 \geq \dots \geq \tilde{\lambda}_k \geq \tilde{\lambda}_{k+1} \geq \dots \geq \tilde{\lambda}_p$.

Moreover, since $(j = 1, \dots, p)$

$$\mathbf{V}\mathbf{Q}\tilde{\mathbf{c}}_j = \tilde{\lambda}_j \tilde{\mathbf{c}}_j ,$$

premultiplying both members by $\tilde{\mathbf{c}}_j' \mathbf{Q}$, we get

$$\tilde{\mathbf{c}}_j' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{c}}_j = \tilde{\lambda}_j \tilde{\mathbf{c}}_j' \mathbf{Q} \tilde{\mathbf{c}}_j = \tilde{\lambda}_j .$$

On the other hand, as \mathbf{V} is positive definite or positive semi-definite and $\mathbf{Q} = \mathbf{Q}'$, $\mathbf{Q} \mathbf{V} \mathbf{Q}$ is also positive definite or positive semi-definite – with $r(\mathbf{Q} \mathbf{V} \mathbf{Q}) = r(\mathbf{V}) = r(\mathbf{Y})$ – and hence $\tilde{\lambda}_j \geq 0$ ($j = 1, \dots, p$).

Thus – setting

$$\begin{aligned} \tilde{\mathbf{C}}_k &= [\tilde{\mathbf{c}}_1 \cdots \tilde{\mathbf{c}}_k] & , & \quad \tilde{\mathbf{C}}_{p-k} = [\tilde{\mathbf{c}}_{k+1} \cdots \tilde{\mathbf{c}}_p] & , & \quad \tilde{\mathbf{C}}_p = [\tilde{\mathbf{c}}_1 \cdots \tilde{\mathbf{c}}_p] \\ \tilde{\mathbf{D}}_k &= \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_k) & , & \quad \tilde{\mathbf{D}}_{p-k} = \text{diag}(\tilde{\lambda}_{k+1}, \dots, \tilde{\lambda}_p) & , & \quad \tilde{\mathbf{D}}_p = \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_p) \end{aligned}$$

– solutions of the system (*) are provided by $\tilde{\mathbf{C}}_k = [\tilde{\mathbf{c}}_1 \cdots \tilde{\mathbf{c}}_k]$ and $\tilde{\mathbf{L}} = \tilde{\mathbf{D}}_k$.

Summing up, first we have

$$\begin{aligned} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_k &= \tilde{\mathbf{C}}_k \tilde{\mathbf{D}}_k & , & \quad \tilde{\mathbf{C}}_k' \mathbf{Q} \tilde{\mathbf{C}}_k &= \mathbf{I}_k , \\ \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_{p-k} &= \tilde{\mathbf{C}}_{p-k} \tilde{\mathbf{D}}_{p-k} & , & \quad \tilde{\mathbf{C}}_{p-k}' \mathbf{Q} \tilde{\mathbf{C}}_{p-k} &= \mathbf{I}_{p-k} , \\ \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_p &= \tilde{\mathbf{C}}_p \tilde{\mathbf{D}}_p & , & \quad \tilde{\mathbf{C}}_p' \mathbf{Q} \tilde{\mathbf{C}}_p &= \mathbf{I}_p \end{aligned}$$

and

$$\begin{aligned} \tilde{\mathbf{C}}_k' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_k &= \tilde{\mathbf{D}}_k & , & \quad \tilde{\mathbf{C}}_k' \mathbf{Q} \tilde{\mathbf{C}}_k &= \mathbf{I}_k , \\ \tilde{\mathbf{C}}_{p-k}' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_{p-k} &= \tilde{\mathbf{D}}_{p-k} & , & \quad \tilde{\mathbf{C}}_{p-k}' \mathbf{Q} \tilde{\mathbf{C}}_{p-k} &= \mathbf{I}_{p-k} , \\ \tilde{\mathbf{C}}_p' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_p &= \tilde{\mathbf{D}}_p & , & \quad \tilde{\mathbf{C}}_p' \mathbf{Q} \tilde{\mathbf{C}}_p &= \mathbf{I}_p . \end{aligned}$$

Then, we have ($\tilde{\mathbf{P}} = \tilde{\mathbf{C}}_k \tilde{\mathbf{C}}_k' \mathbf{Q}$)

$$\begin{aligned} \mathbf{I}_{\mathbf{g} + \tilde{\mathbf{C}}_k} &= \text{tr} \{ \mathbf{V} \mathbf{Q} \tilde{\mathbf{P}} \} & = & \text{tr} \{ \tilde{\mathbf{C}}_k' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_k \} \\ &= \text{tr} \{ \tilde{\mathbf{D}}_k \} & = & \tilde{\lambda}_1 + \dots + \tilde{\lambda}_k , \\ \mathbf{I}_{\mathbf{g} + \tilde{\mathbf{C}}_k^\perp} &= \text{tr} \{ \mathbf{V} \mathbf{Q} (\mathbf{I} - \tilde{\mathbf{P}}) \} & = & \text{tr} \{ \tilde{\mathbf{C}}_{p-k}' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_{p-k} \} \\ &= \text{tr} \{ \tilde{\mathbf{D}}_{p-k} \} & = & \tilde{\lambda}_{k+1} + \dots + \tilde{\lambda}_p , \\ \mathbf{I}_{\mathbf{g}} &= \text{tr} \{ \mathbf{V} \mathbf{Q} \} & = & \text{tr} \{ \tilde{\mathbf{C}}_p' \mathbf{Q} \mathbf{V} \mathbf{Q} \tilde{\mathbf{C}}_p \} \\ &= \text{tr} \{ \tilde{\mathbf{D}}_p \} & = & \tilde{\lambda}_1 + \dots + \tilde{\lambda}_p . \end{aligned}$$

Moreover, the ratio

$$\text{GQRI} = \frac{\mathbf{I}_{\mathbf{g} + \tilde{\mathbf{C}}_k}}{\mathbf{I}_{\mathbf{g}}} = \frac{\tilde{\lambda}_1 + \dots + \tilde{\lambda}_k}{\tilde{\lambda}_1 + \dots + \tilde{\lambda}_p},$$

which denotes that part of the inertia explained by the linear variety $\mathbf{g} + \tilde{\mathbf{C}}_k$, may be used as an index measuring the global quality of representation of $\mathbf{x}_1 + \dots + \mathbf{x}_n$ on $\mathbf{g} + \tilde{\mathbf{C}}_k$.

The eigenvector $\tilde{\mathbf{c}}_j$ ($j = 1, \dots, p$) is called the j th *principal vector*, while the vectors $\mathbf{Q}\tilde{\mathbf{c}}_j$ and $\tilde{\mathbf{y}}_j = \mathbf{Y}\mathbf{Q}\tilde{\mathbf{c}}_j$ are called, respectively, the j th *principal factor* and the j th *principal component*.

REMARK 5. The solutions $\tilde{\mathbf{C}}_k$ and $\tilde{\mathbf{D}}_k$ of the system (*) are not unique.

All other solutions are obtained by the transformations

$$\tilde{\mathbf{C}}_k \rightarrow \tilde{\mathbf{C}}_k \mathbf{T} \quad , \quad \tilde{\mathbf{D}}_k \rightarrow \mathbf{T}' \tilde{\mathbf{D}}_k \mathbf{T}$$

where \mathbf{T} is an orthogonal matrix of order (k, k) .

In fact, from

$$\mathbf{V}\mathbf{Q}\tilde{\mathbf{C}}_k = \tilde{\mathbf{C}}_k \tilde{\mathbf{D}}_k \quad , \quad \tilde{\mathbf{C}}_k' \mathbf{Q}\tilde{\mathbf{C}}_k = \mathbf{I}_k$$

we get $(\mathbf{T}'\mathbf{T} = \mathbf{I}_k = \mathbf{T}\mathbf{T}')$

$$\mathbf{V}\mathbf{Q}(\tilde{\mathbf{C}}_k \mathbf{T}) = (\tilde{\mathbf{C}}_k \mathbf{T})(\mathbf{T}' \tilde{\mathbf{D}}_k \mathbf{T}) \quad , \quad (\mathbf{T}' \tilde{\mathbf{C}}_k') \mathbf{Q}(\tilde{\mathbf{C}}_k \mathbf{T}) = \mathbf{T}'\mathbf{T} = \mathbf{I}_k$$

and inversely.

Notice that, in this case, $\mathbf{T}' \tilde{\mathbf{D}}_k \mathbf{T}$ is symmetric but not more diagonal.

Moreover,

$$\text{tr}(\mathbf{T}' \tilde{\mathbf{C}}_k' \mathbf{Q}\mathbf{V}\mathbf{Q}\tilde{\mathbf{C}}_k \mathbf{T}) = \text{tr}(\mathbf{T}' \tilde{\mathbf{D}}_k \mathbf{T}) = \text{tr}(\mathbf{T}\mathbf{T}' \tilde{\mathbf{D}}_k) = \text{tr}(\tilde{\mathbf{D}}_k) .$$

The solution $\tilde{\mathbf{C}}_k$ is chosen because, as will become apparent below (Section 5.2), it allows us to build up uncorrelated principal components.

It can also be shown that $\tilde{\mathbf{C}}_k$ may be obtained by a step by step procedure that maximizes the inertia explained by a linear variety of increasing dimension (from 1 to k).

5.2 MAIN PROPERTIES OF PRINCIPAL COMPONENTS

1. $\tilde{\mathbf{y}}_j$ ($j = 1, \dots, p$) is a linear combination of the p vectors $\mathbf{y}_1, \dots, \mathbf{y}_p \in \mathbb{R}^n$ the coefficients of which are represented by the elements of the principal factor $\mathbf{Q}\tilde{\mathbf{c}}_j$.

Moreover,

$$[m_1 \dots m_n] \begin{bmatrix} \tilde{y}_{1j} \\ \vdots \\ \tilde{y}_{nj} \end{bmatrix} = \mathbf{u}'\mathbf{M}\tilde{\mathbf{y}}_j = \mathbf{u}'\mathbf{M}\mathbf{Y}\mathbf{Q}\tilde{\mathbf{c}}_j = 0,$$

namely the (weighted) arithmetic mean of $\tilde{\mathbf{y}}_j$ is zero.

2. Setting

$$\tilde{\mathbf{Y}} = [\tilde{\mathbf{y}}_1 \dots \tilde{\mathbf{y}}_p] = \mathbf{Y}\mathbf{Q}\tilde{\mathbf{C}}_p,$$

it follows

$$\tilde{\mathbf{Y}}'\mathbf{M}\tilde{\mathbf{Y}} = \tilde{\mathbf{C}}_p'\mathbf{Q}\mathbf{V}\mathbf{Q}\tilde{\mathbf{C}}_p = \tilde{\mathbf{D}}_p.$$

This last expression indicates the covariance matrix of principal components. It shows that principal components are uncorrelated (orthogonal) with variances given, respectively, by $\tilde{\lambda}_1, \dots, \tilde{\lambda}_p$.

Notice that

$$r(\tilde{\mathbf{Y}}) = r(\tilde{\mathbf{D}}_p) = r(\mathbf{Y})$$

and

$$\text{tr}(\tilde{\mathbf{Y}}'\mathbf{M}\tilde{\mathbf{Y}}) = \text{tr}(\tilde{\mathbf{D}}_p) = \text{tr}(\mathbf{V}\mathbf{Q}).$$

3. Considering the relation ($h = 1, \dots, r$; $r = r(\tilde{\mathbf{D}}_p) = r(\mathbf{Y})$)

$$\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{Q}\tilde{\mathbf{c}}_h = \tilde{\lambda}_h \tilde{\mathbf{c}}_h$$

and premultiplying both members by $\mathbf{Y}\mathbf{Q}$, we obtain

$$\mathbf{Y}\mathbf{Q}\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{Q}\tilde{\mathbf{c}}_h = \tilde{\lambda}_h \mathbf{Y}\mathbf{Q}\tilde{\mathbf{c}}_h$$

or

$$\mathbf{YQY}'\mathbf{M}\tilde{\mathbf{y}}_h = \tilde{\lambda}_h \tilde{\mathbf{y}}_h .$$

This relation shows that $\tilde{\mathbf{y}}_h$ represents an eigenvector of the matrix $\mathbf{YQY}'\mathbf{M}$, obtained for the eigenvalue $\tilde{\lambda}_h > 0$.

Of course, the eigenvalues different from zero of the matrices $\mathbf{YQY}'\mathbf{M}$ and \mathbf{VQ} are the same with the same multiplicities.

4. Setting ($h = 1, \dots, r$)

$$\tilde{\tilde{\mathbf{y}}}_h = \frac{1}{\sqrt{\tilde{\lambda}_h}} \tilde{\mathbf{y}}_h = \frac{1}{\sqrt{\tilde{\lambda}_h}} \mathbf{YQ}\tilde{\mathbf{c}}_h$$

where $\tilde{\tilde{\mathbf{y}}}_h$ denotes the h th standardized principal component, premultiplying both members by $\mathbf{Y}'\mathbf{M}$, we obtain

$$\tilde{\mathbf{c}}_h = \frac{1}{\sqrt{\tilde{\lambda}_h}} \mathbf{Y}'\mathbf{M}\tilde{\tilde{\mathbf{y}}}_h .$$

These two relations, called *transition formulas*, allow us to pass from $\tilde{\mathbf{c}}_h$ to $\tilde{\tilde{\mathbf{y}}}_h$ and vice versa.

5. Since

$$\mathbf{YQ}\tilde{\mathbf{C}}_p = \tilde{\mathbf{Y}},$$

postmultiplying both members by $\tilde{\mathbf{C}}_p'$ we obtain

$$\mathbf{YQ}\tilde{\mathbf{C}}_p\tilde{\mathbf{C}}_p' = \tilde{\mathbf{Y}}\tilde{\mathbf{C}}_p'$$

from which it follows ($\tilde{\mathbf{C}}_p\tilde{\mathbf{C}}_p' = \mathbf{Q}^{-1}$)

$$\mathbf{Y} = \tilde{\mathbf{Y}}\tilde{\mathbf{C}}_p' = [\tilde{\mathbf{y}}_1 \cdots \tilde{\mathbf{y}}_p] \begin{bmatrix} \tilde{\mathbf{c}}_1' \\ \vdots \\ \tilde{\mathbf{c}}_p' \end{bmatrix} = \sum_j \tilde{\mathbf{y}}_j \tilde{\mathbf{c}}_j' = \sum_h \sqrt{\tilde{\lambda}_h} \tilde{\tilde{\mathbf{y}}}_h \tilde{\mathbf{c}}_h' ,$$

the so-called *reconstitution formula* or *singular value decomposition* of the matrix \mathbf{Y} .

Of course, if the summation is limited to the first $h^* < r$ terms, we obtain

an approximated reconstitution of \mathbf{Y} , namely ($h = 1, \dots, h^*$)

$$\mathbf{Y} \cong \sum_h^{h^*} \tilde{\mathbf{y}}_h \tilde{\mathbf{c}}_h' = \sum_h^{h^*} \sqrt{\tilde{\lambda}_h} \tilde{\mathbf{y}}_h \tilde{\mathbf{c}}_h'.$$

6. The cosine of the angle formed by the vectors \mathbf{y}_j ($j = 1, \dots, p$) and $\tilde{\mathbf{y}}_h$ ($h = 1, \dots, r$) – the linear correlation coefficient r_{jh} – is given by

$$\begin{aligned} \cos(\mathbf{y}_j, \tilde{\mathbf{y}}_h) &= \frac{\mathbf{y}_j' \mathbf{M} \tilde{\mathbf{y}}_h}{\sigma_j \sqrt{\tilde{\lambda}_h}} = \frac{\mathbf{u}_j' \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \tilde{\mathbf{c}}_h}{\sigma_j \sqrt{\tilde{\lambda}_h}} = \frac{\mathbf{u}_j' \mathbf{V} \mathbf{Q} \tilde{\mathbf{c}}_h}{\sigma_j \sqrt{\tilde{\lambda}_h}} \\ &= \frac{\mathbf{u}_j' \tilde{\lambda}_h \tilde{\mathbf{c}}_h}{\sigma_j \sqrt{\tilde{\lambda}_h}} = \frac{\mathbf{u}_j' \sqrt{\tilde{\lambda}_h} \tilde{\mathbf{c}}_h}{\sigma_j} = \frac{\sqrt{\tilde{\lambda}_h} \tilde{\mathbf{c}}_h' \mathbf{u}_j}{\sigma_j} = r_{jh}. \end{aligned}$$

7. The orthogonal projection of \mathbf{y}_j ($j = 1, \dots, p$) on the subspace spanned by the principal component $\tilde{\mathbf{y}}_h$ ($h = 1, \dots, r$) is given by

$$\begin{aligned} \tilde{\mathbf{y}}_h (\tilde{\mathbf{y}}_h' \mathbf{M} \tilde{\mathbf{y}}_h)^{-1} \tilde{\mathbf{y}}_h' \mathbf{M} \mathbf{y}_j &= \tilde{\mathbf{y}}_h \frac{1}{\tilde{\lambda}_h} \tilde{\mathbf{c}}_h' \mathbf{Q} \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{u}_j = \tilde{\mathbf{y}}_h \frac{1}{\tilde{\lambda}_h} \tilde{\mathbf{c}}_h' \mathbf{Q} \mathbf{V} \mathbf{u}_j \\ &= \tilde{\mathbf{y}}_h \frac{1}{\tilde{\lambda}_h} \tilde{\lambda}_h \tilde{\mathbf{c}}_h' \mathbf{u}_j = \tilde{\mathbf{y}}_h \tilde{\mathbf{c}}_h' \mathbf{u}_j = \tilde{\mathbf{y}}_h \frac{\sigma_j}{\sqrt{\tilde{\lambda}_h}} r_{jh}. \end{aligned}$$

8. The orthogonal projection of \mathbf{y}_i ($i = 1, \dots, n$) on the subspace spanned by the principal vector $\tilde{\mathbf{c}}_j$ ($j = 1, \dots, p$) is given by

$$\begin{aligned} \tilde{\mathbf{c}}_j (\tilde{\mathbf{c}}_j' \mathbf{Q} \tilde{\mathbf{c}}_j)^{-1} \tilde{\mathbf{c}}_j' \mathbf{Q} \mathbf{y}_i &= \tilde{\mathbf{c}}_j \tilde{\mathbf{c}}_j' \mathbf{Q} \mathbf{Y}' \mathbf{u}_i \\ &= \tilde{\mathbf{c}}_j \tilde{\mathbf{y}}_j' \mathbf{u}_i = \tilde{\mathbf{c}}_j \tilde{\mathbf{y}}_{ij}. \end{aligned}$$

5.3 CHOICE OF THE EUCLIDEAN METRIC IN THE INDIVIDUAL SPACE

The choice of the Euclidean metric in the individual space, the matrix \mathbf{Q} , is probably one of the most delicate problem in PCA.

As we have said above (Section 2.2.1), this choice generally depends on the measurement units and/or the variances of the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$.

First, suppose that the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ are expressed in the same measurement unit and present approximatively the same variance.

In this case, the metric is usually chosen as $\mathbf{Q} = \mathbf{I}_p$, which is equivalent

to perform a PCA on the basis of the covariance matrix \mathbf{V} .

Notice that, if we represent a change in the measurement unit by a constant $s > 0$, we get

$$\begin{aligned} [\mathbf{x}_1 \cdots \mathbf{x}_p] = \mathbf{X} &\rightarrow [s \mathbf{x}_1 \cdots s \mathbf{x}_p] = \mathbf{X}s \\ [\mathbf{y}_1 \cdots \mathbf{y}_p] = \mathbf{Y} &\rightarrow [s \mathbf{y}_1 \cdots s \mathbf{y}_p] = \mathbf{Y}s \\ \mathbf{V}\tilde{\mathbf{c}} = \tilde{\lambda}\tilde{\mathbf{c}} &\rightarrow (s^2\mathbf{V})\tilde{\mathbf{c}} = (s^2\tilde{\lambda})\tilde{\mathbf{c}} \end{aligned}$$

and, thus,

$$\begin{aligned} [\tilde{\mathbf{y}}_1 \cdots \tilde{\mathbf{y}}_p] = \tilde{\mathbf{Y}} &\rightarrow [s \tilde{\mathbf{y}}_1 \cdots s \tilde{\mathbf{y}}_p] = \tilde{\mathbf{Y}}s \\ \tilde{\mathbf{D}}_p &\rightarrow \tilde{\mathbf{D}}_p s^2. \end{aligned}$$

In other words, each new principal component is s times the corresponding old principal component and has a variance s^2 times the corresponding old variance.

Second, suppose again that the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ are expressed in the same measurement unit but present considerably different variances.

In this case, besides the effects of a change in the measurement unit mentioned above, if we perform a PCA on the basis of the covariance matrix \mathbf{V} , those variables whose variances are largest tend to dominate the first few principal components.

To illustrate the point in the simplest way, suppose that we have two variables whose covariance matrix is

$$\mathbf{V} = \begin{bmatrix} 9 & 0.5 \\ 0.5 & 1 \end{bmatrix}.$$

Simple calculations show that $\tilde{\mathbf{y}}_1 = \mathbf{y}_1 0.9981 + \mathbf{y}_2 0.0621$, namely that the first principal component is almost identified with the first variable, the variable with the largest variance.

A solution may be found in choosing the metric

$$\mathbf{Q} = \mathbf{Q}_{1/\sigma^2} = \text{diag} \left(\frac{1}{\sigma_1^2}, \dots, \frac{1}{\sigma_p^2} \right)$$

which is equivalent to standardizing the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ and to perfor-

ming the analysis on the correlation matrix \mathbf{R} .

Yet, results of PCA based on \mathbf{R} are generally different from the corresponding results based on \mathbf{V} .

Third, suppose that the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ are expressed in different measurement units.

In this case, it does not make sense to perform a PCA on the basis of the covariance matrix \mathbf{V} , because operations involving the trace of that matrix have no meaning.

Moreover, setting

$$\mathbf{S} = \text{diag}(s_1, \dots, s_p)$$

where the constants $s_1 > 0, \dots, s_p > 0$ represent a change in the measurement units, we get

$$\begin{aligned} \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_p \end{bmatrix} &= \mathbf{X} &\rightarrow & \begin{bmatrix} s_1 \mathbf{x}_1 & \cdots & s_p \mathbf{x}_p \end{bmatrix} = \mathbf{XS} \\ \begin{bmatrix} \mathbf{y}_1 & \cdots & \mathbf{y}_p \end{bmatrix} &= \mathbf{Y} &\rightarrow & \begin{bmatrix} s_1 \mathbf{y}_1 & \cdots & s_p \mathbf{y}_p \end{bmatrix} = \mathbf{YS} \\ \mathbf{V}\tilde{\mathbf{c}} &= \tilde{\lambda}\tilde{\mathbf{c}} &\rightarrow & (\mathbf{S}\mathbf{V}\mathbf{S})\tilde{\mathbf{c}}^* = \tilde{\lambda}^*\tilde{\mathbf{c}}^* \end{aligned}$$

and, generally, $\tilde{\lambda}^* \neq \tilde{\lambda}$ and $\tilde{\mathbf{c}}^* \neq \tilde{\mathbf{c}}$.

Again, a solution may be found in standardizing the variables $\mathbf{x}_1, \dots, \mathbf{x}_p$ and in performing the analysis on the correlation matrix \mathbf{R} .

6 GRAPHICAL REPRESENTATION OF INDIVIDUALS AND VARIABLES

6.1 GRAPHICAL REPRESENTATION OF INDIVIDUALS

Assuming that $r = r(\mathbf{Y}) \geq 2$, a graphical representation of the n individuals $\mathbf{y}_1, \dots, \mathbf{y}_n$ (measured in terms of deviations from the means) is usually obtained by their orthogonal projections on the subspace \tilde{C}_2 spanned by the first two principal vectors $\tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2$ (*principal plane*).

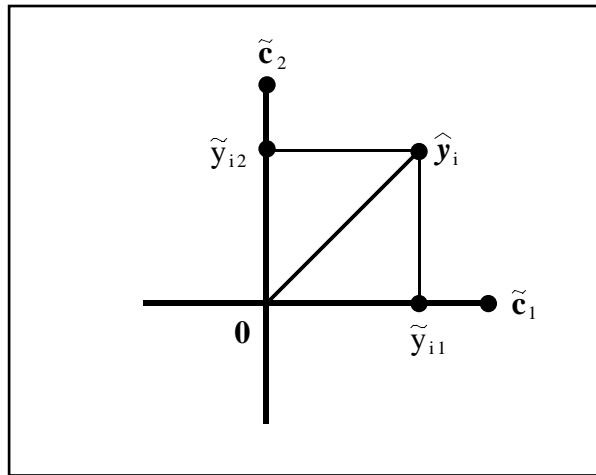
Taking into account what was mentioned above (Section 5.2.8) and denoting by $\hat{\mathbf{y}}_i$ the orthogonal projection of \mathbf{y}_i ($i = 1, \dots, n$) on the principal plane, we have

$$\hat{\mathbf{y}}_i = \tilde{\mathbf{c}}_1 \tilde{y}_{i1} + \tilde{\mathbf{c}}_2 \tilde{y}_{i2}$$

where \tilde{y}_{ij} ($j = 1, 2$) denotes the i th element of the principal component $\tilde{\mathbf{y}}_j$.

Thus, the co-ordinates of $\hat{\mathbf{y}}_i$ relative to $\tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2$ are $(\tilde{y}_{i1}, \tilde{y}_{i2})$ (Fig. 7).

Fig. 7



A measure of the global quality of representation of $\mathbf{y}_1, \dots, \mathbf{y}_n$ on the principal plane is given by the index

$$\text{GQRI} = \frac{\tilde{\lambda}_1 + \tilde{\lambda}_2}{\tilde{\lambda}_1 + \dots + \tilde{\lambda}_p}$$

which may be interpreted as that part of inertia of y_1, \dots, y_n explained by the subspace \tilde{C}_2 (Remark 4).

REMARK 6. Generally, the representation of the individuals on the principal plane is judged to be adequate if GQRI equals or exceeds a predetermined threshold (for example, 0.7). ■

However, as the representation of y_1, \dots, y_n may be good even if some individual y_i is far from its orthogonal projection \hat{y}_i , it is necessary to consider the quality of representation of each y_i ($i = 1, \dots, n$)

An index which may serve this purpose is given by the square cosine of the angle formed by y_i and \hat{y}_i , that is to say by

$$QR(i; \tilde{c}_1, \tilde{c}_2) = \frac{(y_i' Q \hat{y}_i)^2}{(y_i' Q y_i) (\hat{y}_i' Q \hat{y}_i)}.$$

A high $QR(i; \tilde{c}_1, \tilde{c}_2)$ – for example, $QR(i; \tilde{c}_1, \tilde{c}_2) \geq 0.7$ – means that y_i is well represented by \hat{y}_i ; on the contrary, a low $QR(i; \tilde{c}_1, \tilde{c}_2)$ means that the representation of y_i by \hat{y}_i is poor.

Notice that an explicit expression of $QR(i; \tilde{c}_1, \tilde{c}_2)$ may be obtained taking into account that we have the following identities

$$y_i' Q \hat{y}_i = y_i' Q P \tilde{c}_2 y_i = y_i' Q P \tilde{c}_2 P \tilde{c}_2 y_i = y_i' P_{\tilde{c}_2}' Q P \tilde{c}_2 y_i = \hat{y}_i' Q \hat{y}_i$$

where $P_{\tilde{c}_2}$ ($Q P_{\tilde{c}_2} = P_{\tilde{c}_2}' Q$) denotes the orthogonal projection matrix on the subspace \tilde{C}_2 , and

$$\begin{aligned} \hat{y}_i' Q \hat{y}_i &= (\tilde{c}_1 \tilde{y}_{i1} + \tilde{c}_2 \tilde{y}_{i2})' Q (\tilde{c}_1 \tilde{y}_{i1} + \tilde{c}_2 \tilde{y}_{i2}) = \tilde{y}_{i1}^2 + \tilde{y}_{i2}^2, \\ y_i' Q y_i &= u_i' Y Q Y' u_i = u_i' \tilde{Y} \tilde{C}_p' Q \tilde{C}_p \tilde{Y}' u_i = u_i' \tilde{Y} \tilde{Y}' u_i \\ &= [\tilde{y}_{i1} \cdots \tilde{y}_{ip}] \begin{bmatrix} \tilde{y}_{i1} \\ \vdots \\ \tilde{y}_{ip} \end{bmatrix} = \tilde{y}_{i1}^2 + \dots + \tilde{y}_{ip}^2. \end{aligned}$$

Thus,

$$QR(i; \tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2) = \frac{(\mathbf{y}'_i \mathbf{Q} \hat{\mathbf{y}}_i)^2}{(\mathbf{y}'_i \mathbf{Q} \mathbf{y}_i)(\hat{\mathbf{y}}'_i \mathbf{Q} \hat{\mathbf{y}}_i)} = \frac{\hat{\mathbf{y}}'_i \mathbf{Q} \hat{\mathbf{y}}_i}{\mathbf{y}'_i \mathbf{Q} \mathbf{y}_i} = \frac{\tilde{y}_{i1}^2 + \tilde{y}_{i2}^2}{\tilde{y}_{i1}^2 + \dots + \tilde{y}_{ip}^2}.$$

Moreover, since we can write

$$QR(i; \tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2) = \frac{\tilde{y}_{i1}^2}{\mathbf{y}'_i \mathbf{Q} \mathbf{y}_i} + \frac{\tilde{y}_{i2}^2}{\mathbf{y}'_i \mathbf{Q} \mathbf{y}_i} = QR(i; \tilde{\mathbf{c}}_1) + QR(i; \tilde{\mathbf{c}}_2)$$

where $QR(i; \tilde{\mathbf{c}}_j)$ ($j = 1, 2$) denotes the square cosine of the angle formed by \mathbf{y}_i and its orthogonal projection on the subspace spanned by $\tilde{\mathbf{c}}_j$ and is a measure of the quality of representation of \mathbf{y}_i on that subspace, we are able to attribute to each axis the due part of $QR(i; \tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2)$.

After having examined the quality of representation of each individual by means of the index $QR(i; \tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2)$, we are in a position to correctly judge proximities among their orthogonal projections on the principal plane: if two individuals $\mathbf{y}_i, \mathbf{y}_{i^*}$ are close means that $\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_{i^*}$ are close too, provided they are well represented.

In interpreting results of the analysis, it is also important to examine the contribution of each individual \mathbf{y}_i to the inertia $\tilde{\lambda}_j$ explained by $\tilde{\mathbf{c}}_j$.

Since

$$\tilde{\lambda}_j = \tilde{\mathbf{c}}'_j \mathbf{Q} \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \tilde{\mathbf{c}}_j = \tilde{\mathbf{y}}'_j \mathbf{M} \tilde{\mathbf{y}}_j = \sum_i m_i \tilde{y}_{ij}^2,$$

an index often considered is

$$C(i; \tilde{\mathbf{c}}_j) = \frac{m_i \tilde{y}_{ij}^2}{\tilde{\lambda}_j}.$$

The usefulness of examining these contributions may be pointed out first noting that, on the graph, only the co-ordinates of $\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_n$ relative to $\tilde{\mathbf{c}}_j$ are represented; thus, our attention falls on points with a high $|\tilde{y}_{i1}|$, even if some of these may have a small weight.

On the contrary, taking up the examination of contributions – as $C(i; \tilde{\mathbf{c}}_j)$ depends both on m_i and \tilde{y}_{i1}^2 – allows us to detect those individuals which have contributed most to the inertia explained by the subspace under con-

sideration, namely the individuals characterizing that subspace.

Of course, if the individuals all have the same importance – namely, if $\mathbf{M} = \text{diag}(\frac{1}{n}, \dots, \frac{1}{n})$ – examination of the co-ordinates relative to $\tilde{\mathbf{c}}_j$ suffices.

Moreover, it may happen that the contribution of an individual relative to the others is very high; in that case, it is advisable to perform the analysis again after its exclusion from the data set and to reintroduce it as a «supplementary» individual (Section 6.3). This allows us to appreciate differences among remaining individuals, differences which might otherwise be difficult to visualize on the graph since the point scatter is strongly conditioned by the presence of an atypical individual.

6.2 GRAPHICAL REPRESENTATION OF VARIABLES

Assuming that $r = r(\mathbf{Y}) \geq 2$, a graphical representation of the p variables $\mathbf{y}_1, \dots, \mathbf{y}_p$ (measured in terms of deviations from the means) is usually obtained by their orthogonal projections on the subspace $S(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ spanned by the first two standardized principal components $\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2$.

Taking into account what we said above (Section 5.2.7) and denoting by $\hat{\mathbf{y}}_j$ the orthogonal projection of \mathbf{y}_j ($j = 1, \dots, p$) on $S(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$, we have

$$\hat{\mathbf{y}}_j = \tilde{\mathbf{y}}_1 \sigma_j r_{j1} + \tilde{\mathbf{y}}_2 \sigma_j r_{j2}$$

where r_{j1} and r_{j2} denote, respectively, the linear correlation coefficients of \mathbf{y}_j with $\tilde{\mathbf{y}}_1$ and $\tilde{\mathbf{y}}_2$.

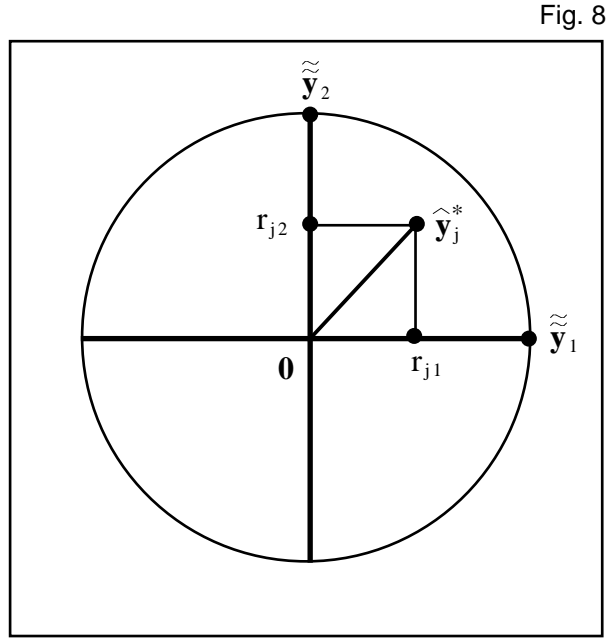
Thus, the co-ordinates of $\hat{\mathbf{y}}_j$ relative to $\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2$ are $(\sigma_j r_{j1}, \sigma_j r_{j2})$.

However, since we are mainly interested in representing linear correlations between pairs of variables or between a variable and a principal component and linear correlations are invariant if each variable is scaled by its standard deviation, it is more suitable to work with standardized variables.

In that case, the orthogonal projection $\hat{\mathbf{y}}_j^*$ of the standardized variable $\mathbf{y}_j^* = \mathbf{y}_j / \sigma_j$ ($j = 1, \dots, p$) on $S(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ is given by

$$\hat{\mathbf{y}}_j^* = \tilde{\mathbf{y}}_1 r_{j1} + \tilde{\mathbf{y}}_2 r_{j2}$$

so that the co-ordinates of $\hat{\mathbf{y}}_j^*$ relative to $\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2$ are (r_{j1}, r_{j2}) (Fig. 8) and hence it is very easy to distinguish those variables which are the most correlated with a principal component and which play a significant role in its interpretation.



Of course, each $\hat{\mathbf{y}}_j^*$ ($j = 1, \dots, p$) is inside a circle of centre $\mathbf{0}$ and radius 1 (the so-called *correlation circle*).

Moreover, the quality of representation of each variable on $S(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ can be judged by means of the square cosine of the angle formed by \mathbf{y}_j^* and $\hat{\mathbf{y}}_j^*$ which is given by $((\mathbf{y}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*)) = 1$

$$QR(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) = \frac{[(\mathbf{y}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*)]^2}{[(\mathbf{y}_j^*)' \mathbf{M}(\mathbf{y}_j^*)][(\hat{\mathbf{y}}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*)]} = \frac{[(\mathbf{y}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*)]^2}{(\hat{\mathbf{y}}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*)}.$$

A high $QR(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ – for example, $QR(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) \geq 0.7$ – means that \mathbf{y}_j^* is well represented by $\hat{\mathbf{y}}_j^*$; on the contrary, a low $QR(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ means that the representation of \mathbf{y}_j^* by $\hat{\mathbf{y}}_j^*$ is poor.

Notice that another expression of $QR(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ may be obtained taking

into account that

$$(\hat{\mathbf{y}}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*) = (\tilde{\mathbf{y}}_1 r_{j1} + \tilde{\mathbf{y}}_2 r_{j2})' \mathbf{M}(\tilde{\mathbf{y}}_1 r_{j1} + \tilde{\mathbf{y}}_2 r_{j2}) = r_{j1}^2 + r_{j2}^2$$

and (Section 5.2.6)

$$\begin{aligned} (\mathbf{y}_j^*)' \mathbf{M}(\hat{\mathbf{y}}_j^*) &= (\mathbf{y}_j^*)' \mathbf{M}(\tilde{\mathbf{y}}_1 r_{j1} + \tilde{\mathbf{y}}_2 r_{j2}) = \frac{\mathbf{y}_j'}{\sigma_j} \mathbf{M} \left(\frac{\tilde{\mathbf{y}}_1}{\sqrt{\lambda_1}} r_{j1} + \frac{\tilde{\mathbf{y}}_2}{\sqrt{\lambda_2}} r_{j2} \right) \\ &= \frac{\mathbf{y}_j' \mathbf{M} \tilde{\mathbf{y}}_1}{\sigma_j \sqrt{\lambda_1}} r_{j1} + \frac{\mathbf{y}_j' \mathbf{M} \tilde{\mathbf{y}}_2}{\sigma_j \sqrt{\lambda_2}} r_{j2} = r_{j1}^2 + r_{j2}^2. \end{aligned}$$

Thus,

$$\text{QR}(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) = r_{j1}^2 + r_{j2}^2.$$

On the other hand, since $\text{QR}(j; \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ also denotes the square distance of $\hat{\mathbf{y}}_j^*$ from the correlation circle centre, we can see that well-represented points lie near the circumference of the correlation circle.

Concluding, for well-represented variables we can visualize on the correlation circle:

- which variables are correlated among themselves and with each principal component;
- which variables are uncorrelated (orthogonal) among themselves and with each principal component.

6.3 SUPPLEMENTARY INDIVIDUALS AND VARIABLES

In applying PCA, it often happens that additional information is available besides that contained in the data matrix \mathbf{Y} .

For example, we may have m additional individuals (measured in terms of deviations from the means)

$$\mathbf{y}'_{n+1} = [y_{n+1,1} \cdots y_{n+1,p}] \quad , \dots \quad , \quad \mathbf{y}'_{n+m} = [y_{n+m,1} \cdots y_{n+m,p}]$$

which belong to a control group and which cannot therefore be included in \mathbf{Y} .

Analogously, we might have q additional variables (measured in terms of deviations from the means)

$$\mathbf{y}'_{p+1} = [y_{1,p+1} \cdots y_{n,p+1}] , \dots , \mathbf{y}'_{p+q} = [y_{1,p+q} \cdots y_{n,p+q}]$$

which are of a different nature with respect to the variables contained in \mathbf{Y} and which we do not wish to incorporate in \mathbf{Y} .

After having obtained – on the basis of the matrix \mathbf{Y} – principal vectors and principal components and represented the initial individuals and variables, we would like to place the m additional individuals and the q additional variables on the respective graphs.

The procedure for doing this consists of positioning the supplementary individuals and variables on the graphs.

Of course, as before, the orthogonal projection matrix of a supplementary individual on the subspace spanned by $\tilde{\mathbf{c}}_1, \tilde{\mathbf{c}}_2$ is

$$[\tilde{\mathbf{c}}_1 \ \tilde{\mathbf{c}}_2][\tilde{\mathbf{c}}_1 \ \tilde{\mathbf{c}}_2]' \mathbf{Q} ,$$

and the orthogonal projection matrix of a supplementary variable on the subspace spanned by the first two standardized principal components $\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2$ is

$$[\tilde{\mathbf{y}}_1 \ \tilde{\mathbf{y}}_2][\tilde{\mathbf{y}}_1 \ \tilde{\mathbf{y}}_2]' \mathbf{M} .$$

7 OTHER APPROACHES TO PCA

7.1 THE APPROACH IN TERMS OF WEIGHTED SUM OF SQUARE DISTANCES BETWEEN ANY PAIR OF INDIVIDUALS

Consider the n individuals y_1, \dots, y_n (measured in terms of deviations from the means) with weights given, respectively, by m_1, \dots, m_n .

Let

$$d_{i,i^*}^2 = (y_i - y_{i^*})' Q (y_i - y_{i^*})$$

the square distance between a pair of individuals (y_i, y_{i^*}) ($i, i^* = 1, \dots, n$) and

$$D^2 = \begin{bmatrix} d_{1,1}^2 & \dots & d_{1,n}^2 \\ \dots & \dots & \dots \\ d_{n,1}^2 & \dots & d_{n,n}^2 \end{bmatrix}$$

the corresponding square distance matrix ⁽⁶⁾.

Setting

$$d^2 u' = \begin{bmatrix} y_1' Q y_1 \\ \vdots \\ y_n' Q y_n \end{bmatrix} u' = \begin{bmatrix} y_1' Q y_1 & \dots & y_1' Q y_n \\ \dots & \dots & \dots \\ y_n' Q y_1 & \dots & y_n' Q y_n \end{bmatrix},$$

$$u d^{2'} = u [y_1' Q y_1 \dots y_n' Q y_n] = \begin{bmatrix} y_1' Q y_1 & \dots & y_n' Q y_n \\ \dots & \dots & \dots \\ y_1' Q y_1 & \dots & y_n' Q y_n \end{bmatrix},$$

$$Y Q Y' = \begin{bmatrix} y_1' \\ \vdots \\ y_n' \end{bmatrix} Q [y_1 \dots y_n] = \begin{bmatrix} y_1' Q y_1 & \dots & y_1' Q y_n \\ \dots & \dots & \dots \\ y_n' Q y_1 & \dots & y_n' Q y_n \end{bmatrix},$$

it can easily be shown that

(6) Of course, this matrix is symmetric and has the elements on the principal diagonal all equal to zero.

$$\mathbf{D}^2 = \mathbf{d}^2 \mathbf{u}' + \mathbf{u} \mathbf{d}^{2'} - 2 \mathbf{Y} \mathbf{Q} \mathbf{Y}' .$$

Now, consider the weighted sum of square distances between any pair of individuals, each square distance weighted by $m_i m_{i^*}$.

We want to show that this last quantity, which can be written as $\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{D}^2)$, is nothing other than twice \mathbf{I}_g .

In fact, first we have ($\mathbf{u}' \mathbf{M} \mathbf{Y} = \mathbf{0}_{(1,p)}$)

$$\begin{aligned} \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{D}^2) &= \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} (\mathbf{d}^2 \mathbf{u}' + \mathbf{u} \mathbf{d}^{2'} - 2 \mathbf{Y} \mathbf{Q} \mathbf{Y}')) \\ &= \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{d}^2 \mathbf{u}') + \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{u} \mathbf{d}^{2'}) - 2 \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{Y} \mathbf{Q} \mathbf{Y}') \\ &= \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{d}^2 \mathbf{u}') + \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{u} \mathbf{d}^{2'}) . \end{aligned}$$

Then – since

$$\begin{aligned} \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{d}^2 \mathbf{u}') &= \text{tr}(\mathbf{d}^2 \mathbf{u}' \mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M}) \\ &= \text{tr}(\mathbf{d}^2 \mathbf{u}' \mathbf{M}) = \text{tr} \begin{bmatrix} m_1 \mathbf{y}'_1 \mathbf{Q} \mathbf{y}_1 \cdots m_n \mathbf{y}'_1 \mathbf{Q} \mathbf{y}_1 \\ \cdots \cdots \cdots \\ m_1 \mathbf{y}'_n \mathbf{Q} \mathbf{y}_n \cdots m_n \mathbf{y}'_n \mathbf{Q} \mathbf{y}_n \end{bmatrix} \\ &= \text{tr} \begin{bmatrix} m_1 \mathbf{y}'_1 \mathbf{Q} \mathbf{y}_1 \cdots m_1 \mathbf{y}'_1 \mathbf{Q} \mathbf{y}_n \\ \cdots \cdots \cdots \\ m_n \mathbf{y}'_n \mathbf{Q} \mathbf{y}_1 \cdots m_n \mathbf{y}'_n \mathbf{Q} \mathbf{y}_n \end{bmatrix} = \text{tr}(\mathbf{Y} \mathbf{Q} \mathbf{Y}' \mathbf{M}) \\ &= \text{tr}(\mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q}) = \text{tr}(\mathbf{V} \mathbf{Q}) \\ &= \mathbf{I}_g \end{aligned}$$

and, analogously,

$$\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{u} \mathbf{d}^{2'}) = \mathbf{I}_g ,$$

– we get

$$\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{D}^2) = 2 \mathbf{I}_g .$$

Next, let

$$\widehat{d}_{i,i^*}^2 = (\mathbf{P} \mathbf{y}_i - \mathbf{P} \mathbf{y}_{i^*})' \mathbf{Q} (\mathbf{P} \mathbf{y}_i - \mathbf{P} \mathbf{y}_{i^*})$$

the square distance between a pair of projected individuals $(\mathbf{P} \mathbf{y}_i, \mathbf{P} \mathbf{y}_{i^*})$ on \mathbf{C}_k and

$$\widehat{\mathbf{D}}^2 = \begin{bmatrix} \widehat{d}_{1,1}^2 & \cdots & \widehat{d}_{1,n}^2 \\ \cdots & \cdots & \cdots \\ \widehat{d}_{n,1}^2 & \cdots & \widehat{d}_{n,n}^2 \end{bmatrix}$$

the corresponding square distance matrix.

Setting

$$\widehat{\mathbf{d}}^2 \mathbf{u}' = \begin{bmatrix} \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \end{bmatrix} \mathbf{u}' = \begin{bmatrix} \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \cdots \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \\ \cdots \quad \cdots \quad \cdots \\ \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \cdots \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \end{bmatrix},$$

$$\mathbf{u} \widehat{\mathbf{d}}^{2'} = \mathbf{u} [\mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \cdots \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n] = \begin{bmatrix} \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \cdots \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \\ \cdots \quad \cdots \quad \cdots \\ \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \cdots \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \end{bmatrix},$$

$$\mathbf{Y} \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{Y}' = \begin{bmatrix} \mathbf{y}'_1 \mathbf{P}' \\ \vdots \\ \mathbf{y}'_n \mathbf{P}' \end{bmatrix} \mathbf{Q} [\mathbf{P} \mathbf{y}_1 \cdots \mathbf{P} \mathbf{y}_n] = \begin{bmatrix} \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \cdots \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \\ \cdots \quad \cdots \quad \cdots \\ \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \cdots \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \end{bmatrix},$$

it can easily be shown that

$$\widehat{\mathbf{D}}^2 = \widehat{\mathbf{d}}^2 \mathbf{u}' + \mathbf{u} \widehat{\mathbf{d}}^{2'} - 2 \mathbf{Y} \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{Y}' .$$

Now, consider the weighted sum of square distances between any pair of projected individuals, each square distance weighted by $m_i m_{i^*}$.

We want to show that this last quantity, which can be written as $\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \widehat{\mathbf{D}}^2)$, is nothing other than twice $\mathbf{I}_{\mathbf{g} + \mathbf{C}_k}$.

In fact, first we have ($\mathbf{u}' \mathbf{M} \mathbf{Y} = \mathbf{0}_{(1,p)}$)

$$\begin{aligned} \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \widehat{\mathbf{D}}^2) &= \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} (\widehat{\mathbf{d}}^2 \mathbf{u}' + \mathbf{u} \widehat{\mathbf{d}}^{2'} - 2 \mathbf{Y} \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{Y}')) \\ &= \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \widehat{\mathbf{d}}^2 \mathbf{u}') + \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{u} \widehat{\mathbf{d}}^{2'}) - 2 \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{Y} \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{Y}') \\ &= \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \widehat{\mathbf{d}}^2 \mathbf{u}') + \text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \mathbf{u} \widehat{\mathbf{d}}^{2'}) . \end{aligned}$$

Then – since

$$\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \widehat{\mathbf{d}}^2 \mathbf{u}') = \text{tr}(\widehat{\mathbf{d}}^2 \mathbf{u}' \mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M})$$

$$\begin{aligned}
&= \text{tr}(\widehat{\mathbf{d}}^2 \mathbf{u}'\mathbf{M}) &= \text{tr} \begin{bmatrix} m_1 \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 & \cdots & m_n \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 \\ \cdots & \cdots & \cdots \\ m_1 \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n & \cdots & m_n \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \end{bmatrix} \\
&= \text{tr} \begin{bmatrix} m_1 \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 & \cdots & m_1 \mathbf{y}'_1 \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \\ \cdots & \cdots & \cdots \\ m_n \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_1 & \cdots & m_n \mathbf{y}'_n \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{y}_n \end{bmatrix} &= \text{tr}(\mathbf{Y} \mathbf{P}' \mathbf{Q} \mathbf{P} \mathbf{Y}' \mathbf{M}) \\
&= \text{tr}(\mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \mathbf{P}) &= \text{tr}(\mathbf{V} \mathbf{Q} \mathbf{P}) \\
&= \mathbf{I}_{\mathbf{g} + \mathbf{C}_k}
\end{aligned}$$

and, analogously,

$$\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \widehat{\mathbf{d}}^2) = \mathbf{I}_{\mathbf{g} + \mathbf{C}_k},$$

– we get

$$\text{tr}(\mathbf{M} \mathbf{u} \mathbf{u}' \mathbf{M} \widehat{\mathbf{D}}^2) = 2\mathbf{I}_{\mathbf{g} + \mathbf{C}_k}.$$

Summing up, since we have shown that the weighted sum of square distances between any pair of individuals is twice $\mathbf{I}_{\mathbf{g}}$ and that the weighted sum of square distances between any pair of projected individuals is twice $\mathbf{I}_{\mathbf{g} + \mathbf{C}_k}$, finding out the subspace maximizing the inertia explained by $\mathbf{I}_{\mathbf{g} + \mathbf{C}_k}$ (Section 5.1) is equivalent to looking for the subspace maximizing the weighted sum of square distances between any pair of projected individuals.

In other words, with this interpretation the criterion consists of finding out the subspace modifying as little as possible the weighted sum of square distances between any pair of individuals when passing to \mathbf{C}_k .

7.2 THE APPROACH IN TERMS OF GLOBAL VARIABILITY AND GENERALIZED VARIANCE

Consider the problem of finding out (Section 5.1)

$$\text{Max}_{\mathbf{C}_k} \text{tr} \{ \mathbf{C}'_k \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \} = \text{Max}_{\mathbf{C}_k} \text{tr} \{ \mathbf{C}'_k \mathbf{Q} \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \mathbf{C}_k \}, \quad \mathbf{C}'_k \mathbf{Q} \mathbf{C}_k = \mathbf{I}_k.$$

Since $\mathbf{c}_h' \mathbf{Q} \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \mathbf{c}_h$ is the variance of the h th ($h = 1, \dots, k$) linear combination $\mathbf{Y} \mathbf{Q} \mathbf{c}_h$ of the p variables $\mathbf{y}_1, \dots, \mathbf{y}_p$ (measured in terms of deviations from the means), we realize that the above-mentioned problem is equivalent to maximizing the *global variability* of the k linear combinations $\mathbf{Y} \mathbf{Q} \mathbf{c}_1, \dots, \mathbf{Y} \mathbf{Q} \mathbf{c}_k$ of $\mathbf{y}_1, \dots, \mathbf{y}_p$ under the constraint $\mathbf{C}'_k \mathbf{Q} \mathbf{C}_k = \mathbf{I}_k$.

Of course, $\tilde{\mathbf{C}}_k = [\tilde{\mathbf{c}}_1 \dots \tilde{\mathbf{c}}_k]$ is a solution of this problem and

$$\text{tr} \{ \tilde{\mathbf{C}}_k' \mathbf{V} \tilde{\mathbf{C}}_k \} = \sum_h \tilde{\lambda}_h.$$

Instead of considering $\text{tr} \{ \mathbf{C}'_k \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \}$ as a measure of the variability of the k linear combinations $\mathbf{Y} \mathbf{Q} \mathbf{c}_1, \dots, \mathbf{Y} \mathbf{Q} \mathbf{c}_k$ of $\mathbf{y}_1, \dots, \mathbf{y}_p$, we also may refer to $\det \{ \mathbf{C}'_k \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \}$, namely to the so-called *generalized variance*.

In this case, the problem becomes

$$\text{Max}_{\mathbf{C}_k} \det \{ \mathbf{C}'_k \mathbf{Q} \mathbf{V} \mathbf{Q} \mathbf{C}_k \}, \quad \mathbf{C}'_k \mathbf{Q} \mathbf{C}_k = \mathbf{I}_k.$$

It can be shown ⁽⁷⁾ that $\tilde{\mathbf{C}}_k = [\tilde{\mathbf{c}}_1 \dots \tilde{\mathbf{c}}_k]$ is a solution and that

$$\det \{ \tilde{\mathbf{C}}_k' \mathbf{V} \tilde{\mathbf{C}}_k \} = \prod_h \tilde{\lambda}_h.$$

7.3 THE APPROACH IN TERMS OF SUM OF SQUARE LINEAR CORRELATION COEFFICIENTS BETWEEN A NORMALIZED LINEAR COMBINATION OF THE ORIGINAL VARIABLES AND EACH ORIGINAL VARIABLE

Consider the problem of finding out a normalized variable $\tilde{\mathbf{y}}_{(1)}$, linear combination of $\mathbf{y}_1, \dots, \mathbf{y}_p$, maximizing the sum of the square linear correlation coefficients between $\tilde{\mathbf{y}}_{(1)}$ and each \mathbf{y}_j ($j = 1, \dots, p$).

Denote by $\tilde{\mathbf{y}}$ a generic normalized linear combination of $\mathbf{y}_1, \dots, \mathbf{y}_p$.

Since we have ($\tilde{\mathbf{y}}' \mathbf{M} \tilde{\mathbf{y}} = 1$)

(7) A proof is given in Jolliffe on pp.15-16.

$$\begin{aligned}
\sum_j \cos^2(\tilde{\mathbf{y}}, \mathbf{y}_j) &= \sum_j \frac{(\tilde{\mathbf{y}}' \mathbf{M} \mathbf{y}_j)^2}{\sigma_j^2} \\
&= \sum_j \frac{\tilde{\mathbf{y}}' \mathbf{M} \mathbf{y}_j \mathbf{y}_j' \mathbf{M} \tilde{\mathbf{y}}}{\sigma_j^2} = \tilde{\mathbf{y}}' \mathbf{M} \left(\sum_j \frac{\mathbf{y}_j \mathbf{y}_j'}{\sigma_j^2} \right) \mathbf{M} \tilde{\mathbf{y}} \\
&= \tilde{\mathbf{y}}' \mathbf{M} \left(\begin{bmatrix} \mathbf{y}_1 & \cdots & \mathbf{y}_p \end{bmatrix} \mathbf{Q}_{1/\sigma^2} \begin{bmatrix} \mathbf{y}_1' \\ \vdots \\ \mathbf{y}_p' \end{bmatrix} \right) \mathbf{M} \tilde{\mathbf{y}} = \tilde{\mathbf{y}}' \mathbf{M} \mathbf{Y} \mathbf{Q}_{1/\sigma^2} \mathbf{Y}' \mathbf{M} \tilde{\mathbf{y}},
\end{aligned}$$

we must find out

$$\text{Max}_{\tilde{\mathbf{y}}} \tilde{\mathbf{y}}' \mathbf{M} \mathbf{Y} \mathbf{Q}_{1/\sigma^2} \mathbf{Y}' \mathbf{M} \tilde{\mathbf{y}}, \quad \tilde{\mathbf{y}}' \mathbf{M} \tilde{\mathbf{y}} = 1.$$

As can easily be seen, a solution of this problem is given by the normalized eigenvector $\tilde{\mathbf{y}}_{(1)}$ of the matrix $\mathbf{Y} \mathbf{Q}_{1/\sigma^2} \mathbf{Y}' \mathbf{M}$ associated with the eigenvalue $\hat{\lambda}_1$, and that $\sum_j \cos^2(\tilde{\mathbf{y}}, \mathbf{y}_j) = \hat{\lambda}_1$.

Thus, assuming that $\mathbf{Q} = \mathbf{Q}_{1/\sigma^2}$, $\tilde{\mathbf{y}}_{(1)}$ equals $\tilde{\mathbf{y}}_1$, the first standardized principal component.

Of course, an analogous meaning may be attributed to each of the subsequent standardized principal components.

7.4 THE APPROACH IN TERMS OF THE MULTIVARIABLE LINEAR MODEL

The other approach we would like to mention is based on the multivariable linear model.

To this end, first remember that this latter can be expressed in the form

$$\mathbf{Y}_1 = \mathbf{Y}_2 \mathbf{H}_2 + \mathbf{E}_2$$

where

- \mathbf{Y}_1 is the matrix, of order (n, p_1) , of the observed values of p_1 dependent variables, measured in terms of deviations from the means;
- \mathbf{Y}_2 is the matrix, of order (n, p_2) , of the observed values of p_2 independent variables, measured in terms of deviations from the means;

- \mathbf{H}_2 is a matrix, of order (p_2, p_1) , of unknown coefficients;
- \mathbf{E}_2 is a matrix, of order (n, p_1) , of «residuals».

In order to determine the matrix \mathbf{H}_2 , we can choose a least squares criterion, which means finding out

$$\text{Min}_{\mathbf{H}_2} \text{tr} \{(\mathbf{Y}_1 - \mathbf{Y}_2 \mathbf{H}_2)' \mathbf{M} (\mathbf{Y}_1 - \mathbf{Y}_2 \mathbf{H}_2) \mathbf{Q}\} .$$

Of course, assuming that $r(\mathbf{Y}_2) = p_2$, the best solution is given by $\widehat{\mathbf{H}}_2 = (\mathbf{Y}_2' \mathbf{M} \mathbf{Y}_2)^{-1} \mathbf{Y}_2' \mathbf{M} \mathbf{Y}_1$.

Now, in case $\mathbf{Y}_1 = \mathbf{Y}_2 = \mathbf{Y}$, consider the model

$$\mathbf{Y} = \mathbf{YH} + \mathbf{E}$$

from which is clear that, without any assumption regarding the matrix \mathbf{H} , of order (p, p) , the best solution is trivially given by $\widetilde{\mathbf{H}} = \mathbf{I}_p$.

Then, assume that \mathbf{H} has rank $h^* < r = r(\mathbf{Y})$, so that it may be written in the form (\mathbf{F} and \mathbf{G} of order, respectively, (p, h^*) and (h^*, p))

$$\mathbf{H} = \mathbf{FG}$$

with $r(\mathbf{F}) = r(\mathbf{G}) = h^*$.

Our model becomes

$$\mathbf{Y} = \mathbf{YFG} + \mathbf{E}$$

and we propose to find out

$$\text{Min}_{\mathbf{F}, \mathbf{G}} \text{tr} \{(\mathbf{Y} - \mathbf{YFG})' \mathbf{M} (\mathbf{Y} - \mathbf{YFG}) \mathbf{Q}\} \quad , \quad \mathbf{F}' \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{F} = \mathbf{I}_{h^*} .$$

To this end, first notice that, taking into account the constraint on the matrix \mathbf{F} , our problem lies in finding out

$$\text{Min}_{\mathbf{F}, \mathbf{G}} \{ \text{tr} \{ \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{Q} \} - 2 \text{tr} \{ \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{F} \mathbf{G} \mathbf{Q} \} + \text{tr} \{ \mathbf{G}' \mathbf{G} \mathbf{Q} \} \} \quad , \quad \mathbf{F}' \mathbf{Y}' \mathbf{M} \mathbf{Y} \mathbf{F} = \mathbf{I}_{h^*}$$

or, equivalently,

$$\text{Max}_{\mathbf{F}, \mathbf{G}} \{2\text{tr}\{\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F}\mathbf{G}\mathbf{Q}\} - \text{tr}\{\mathbf{G}'\mathbf{G}\mathbf{Q}\}\} , \quad \mathbf{F}'\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F} = \mathbf{I}_{h^*} .$$

Now, consider the Lagrange function

$$L(\mathbf{F}, \mathbf{G}, \mathbf{L}) = 2\text{tr}\{\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F}\mathbf{G}\mathbf{Q}\} - \text{tr}\{\mathbf{G}'\mathbf{G}\mathbf{Q}\} - \text{tr}\{(\mathbf{F}'\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F} - \mathbf{I}_{h^*}) \mathbf{L}\}$$

where $\mathbf{L} = \mathbf{L}'$ is a matrix of Lagrange multipliers of order (h^*, h^*) .

At $(\tilde{\mathbf{F}}, \tilde{\mathbf{G}}, \tilde{\mathbf{L}})$ where $L(\mathbf{F}, \mathbf{G}, \mathbf{L})$ has a maximum, as can easily be verified, it must be

$$\begin{aligned} \mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{Q}\tilde{\mathbf{G}}' &= \mathbf{Y}'\mathbf{M}\mathbf{Y}\tilde{\mathbf{F}}\tilde{\mathbf{L}} \\ \tilde{\mathbf{F}}'\mathbf{Y}'\mathbf{M}\mathbf{Y} &= \tilde{\mathbf{G}} \\ \tilde{\mathbf{F}}'\mathbf{Y}'\mathbf{M}\mathbf{Y}\tilde{\mathbf{F}} &= \mathbf{I}_{h^*} . \end{aligned}$$

Therefore, we must find out solutions of the system

$$\begin{aligned} \mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{Q}\mathbf{G}' &= \mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F}\mathbf{L} \\ \mathbf{F}'\mathbf{Y}'\mathbf{M}\mathbf{Y} &= \mathbf{G} \\ \mathbf{F}'\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F} &= \mathbf{I}_{h^*} \end{aligned}$$

in the unknowns $\mathbf{F}, \mathbf{G}, \mathbf{L}$.

But, premultiplying the first equation by \mathbf{F}' and taking into account the remaining equations, we obtain

$$\mathbf{F}'\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{Q}\mathbf{Y}'\mathbf{M}\mathbf{Y}\mathbf{F} = \mathbf{L} .$$

This matrix has n eigenvalues $\tilde{\lambda}_1, \dots, \tilde{\lambda}_n$ of which $r = r(\mathbf{Y})$ are positive, the remainder zero.

Then, associate to the first h^* positive eigenvalues $\tilde{\lambda}_1, \dots, \tilde{\lambda}_{h^*}$ the h^* orthonormal eigenvectors $\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_{h^*}$.

Setting

$$\tilde{\mathbf{D}}_{h^*} = \text{diag}(\tilde{\lambda}_1, \dots, \tilde{\lambda}_{h^*}) , \quad \tilde{\mathbf{Y}}_{h^*} = [\tilde{\mathbf{y}}_1 \cdots \tilde{\mathbf{y}}_{h^*}] ,$$

we can write

$$\mathbf{Y}\mathbf{Q}\mathbf{Y}'\mathbf{M}\tilde{\mathbf{Y}}_{h^*} = \tilde{\mathbf{Y}}_{h^*}'\tilde{\mathbf{D}}_{h^*} , \quad \tilde{\mathbf{Y}}_{h^*}'\mathbf{M}\tilde{\mathbf{Y}}_{h^*} = \mathbf{I}_{h^*}$$

and also

$$\tilde{\mathbf{Y}}_{h^*}' \mathbf{M} \mathbf{Y} \mathbf{Q} \mathbf{Y}' \mathbf{M} \tilde{\mathbf{Y}}_{h^*} = \tilde{\mathbf{D}}_{h^*} .$$

Thus, $\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_{h^*}$ are the first h^* standardized principal components of \mathbf{VQ} . Thus, we immediately realize that

$$\tilde{\mathbf{F}} = \mathbf{Q} \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{D}}_{h^*}^{-1/2} , \quad \tilde{\mathbf{G}} = \tilde{\mathbf{F}}' \mathbf{Y}' \mathbf{M} \mathbf{Y} = \tilde{\mathbf{Y}}_{h^*}' \mathbf{M} \mathbf{Y} , \quad \tilde{\mathbf{L}} = \tilde{\mathbf{D}}_{h^*}$$

represent a solution of our problem.

Summing up, both $\tilde{\mathbf{F}}$ and $\tilde{\mathbf{G}}$ can be interpreted in terms of principal factors and standardized principal components.

In fact, while $\tilde{\mathbf{F}}$ is linked to the matrix $\mathbf{Q} \tilde{\mathbf{C}}_{h^*}$ of the principal factors by means of the matrix $\tilde{\mathbf{D}}_{h^*}^{-1/2}$, $\tilde{\mathbf{G}}$, as is at once apparent, represents the matrix of the coefficients of the orthogonal projection of \mathbf{Y} on the subspace spanned by the column vectors of the matrix

$$\tilde{\mathbf{Y}}_{h^*} = \mathbf{Y} \mathbf{Q} \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{D}}_{h^*}^{-1/2} = \mathbf{Y} \tilde{\mathbf{F}} .$$

Notice that, since $(\mathbf{VQ} \tilde{\mathbf{C}}_{h^*} = \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{D}}_{h^*})$

$$\begin{aligned} \tilde{\mathbf{H}} &= \tilde{\mathbf{F}} \tilde{\mathbf{G}} &&= \mathbf{Q} \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{D}}_{h^*}^{-1/2} \tilde{\mathbf{D}}_{h^*}^{-1/2} \tilde{\mathbf{C}}_{h^*}' \mathbf{Q} \mathbf{V} \\ &= \mathbf{Q} \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{D}}_{h^*}^{-1} \tilde{\mathbf{D}}_{h^*} \tilde{\mathbf{C}}_{h^*}' &&= \mathbf{Q} \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{C}}_{h^*}' , \end{aligned}$$

we get

$$\mathbf{Y} \tilde{\mathbf{H}} = \mathbf{Y} \mathbf{Q} \tilde{\mathbf{C}}_{h^*} \tilde{\mathbf{C}}_{h^*}' = \tilde{\mathbf{Y}}_{h^*} \tilde{\mathbf{C}}_{h^*}' = \sum_{h^*} \tilde{\mathbf{y}}_h \tilde{\mathbf{c}}_h' .$$

Namely, $\mathbf{Y} \tilde{\mathbf{H}}$ is an approximated reconstitution of \mathbf{Y} (Section 5.2.5).

REFERENCES

- [1] Anderson, T.W., *Introduction to Multivariate Statistical Analysis*, John Wiley and Sons, New York, 1958.
- [2] Basilevsky, A., *Statistical Factor Analysis and Related Methods*, John Wiley and Sons, New York, 1994.
- [3] Bertier, P., Bouroche, J.M., *Analyse des données multidimensionnelles*, PUF, Paris, 1977.
- [4] Bolasco, S., *Analisi multidimensionale dei dati*, Carocci, Roma, 1999.
- [5] Bouroche, J.M., Saporta, G., *L'analisi dei dati*, CLU, Napoli, 1983.
- [6] Cailliez, F., Pages, G.P., *Introduction à l'analyse des données*, Smash, Paris, 1976.
- [7] Coppi, R., *Appunti di statistica metodologica: analisi lineare dei dati*, Dipartimento di Statistica, Probabilità e Statistiche Applicate, Roma, 1986.
- [8] De Antoni, F., (a cura di), *I fondamenti dell'analisi dei dati*, Istituto di Statistica e Ricerca Sociale "C. Gini", Roma, 1982.
- [9] Delvecchio, F., *Analisi statistica di dati multidimensionali*, Cacucci Editore, Bari, 1992.
- [10] Diday, E., Lemaire, J., Pouget, J., Testu, F., *Eléments d'analyse des données*, Dunod, Paris, 1982.
- [11] Fabbris, L., *Analisi esplorativa di dati multidimensionali*, cleup editore, Padova, 1990.
- [12] Jackson, J.E., *A User's Guide to Principal Components*, John Wiley and Sons, New York, 1991.
- [13] Jolliffe, I.T., *Principal Component Analysis*, Springer-Verlag, New

- York, 2002.
- [14] Krzanowski, W.J., *Principles of Multivariate Analysis*, Oxford University Press, Oxford, 2000.
- [15] Kshirsagar, A.M., *Multivariate Analysis*, Marcel Dekker, Inc., New York, 1972.
- [16] Lebart, L., Morineau, A., Warwick, K.M., *Multivariate Descriptive Analysis*, John Wiley and Sons, New York, 1984.
- [17] Leoni, R., *Alcuni argomenti di analisi statistica multivariata*, Dipartimento Statistico, Firenze, 1978.
- [18] Leoni, R., (a cura di) *Alcuni lavori di analisi statistica multivariata*, SIS, Firenze, 1982.
- [19] Leoni, R., *Principal Component Analysis*, in «Methods for Multi-dimensional Data Analysis», European Courses in Advanced Statistics, Anacapri, 1987.
- [20] Leoni, R., *Algebra lineare per le applicazioni statistiche*, Dipartimento di Statistica "G. Parenti", Firenze, 2007 (sta in <<http://www.ds.unifi.it>> alla voce *Materiale Didattico*).
- [21] Leoni, R., *Modello lineare multivariato e analisi statistica multidimensionale*, in «Conferenze di statistica nell'anno del 750° anniversario dell'Università degli Studi di Siena», Dipartimento di Metodi Quantitativi, Siena, 1994.
- [22] Marchetti, G., *Analisi in componenti principali e approssimazioni di matrici*, Dipartimento Statistico, Firenze, 1984.
- [23] Mardia, K.V., Kent, I.T., Bibby, J.M., *Multivariate Analysis*, Academic Press, London, 1979.
- [24] Mignani, S., Montanari, A., *Appunti di analisi statistica multivariata*, Società Editrice Esculapio, Bologna, 1998.
- [25] Rao, C.R., *The Use and Interpretation of Principal Component*

Analysis in Applied Research, Sankyā A, 26, 1964.

- [26] Rao, C.R., *Matrix Approximations and Reduction of Dimensionality in Multivariate Statistical Analysis*, in «Multivariate Analysis-V» (Krishnaiah, P.R., ed.), North-Holland Publishing Company, Amsterdam, 1980.
- [27] Rencher, A.C., *Methods of Multivariate Analysis*, John Wiley & Sons, New York, 1995.
- [28] Rizzi, A., *Analisi dei dati*, NIS, Roma, 1985.
- [29] Saporta, G., *Probabilités, Analyse des données et Statistique*, Éditions Technip, Paris, 1990.
- [30] Seber, G.A.F., *Multivariate Observations*, John Wiley & Sons, New York, 1984.
- [31] Volle, M., *Analyse des données*, Economica, Paris, 1981.
- [32] Zani, S., *Analisi dei dati statistici II*, Giuffrè Editore, Milano, 2000.