

Principal Component Analysis with external information on both subjects and variables

Analisi Statistica dei Dati Multidimensionali¹

¹Corso di Laurea in Scienze Statistiche e Attuariali

Facoltà di Scienze Economiche e Aziendali
Università degli Studi del Sannio

- $Im(\mathbf{X})$ denotes the subspace spanned by the column vectors of \mathbf{X} .
- Let $Ker(\mathbf{X}^T)$ be the kernel (null space) of \mathbf{X}^T .
- The data matrix \mathbf{X} is often accompanied by external information concerning the rows and columns of data matrix.
- In order to consider several hypotheses and evaluate different effects, we consider several external information \mathbf{H} and \mathbf{Z} for the units and variables of the matrix \mathbf{X} .
- Let $\mathbf{H}_\mathbf{X}$ be the external $n \times c$ row (units) information matrices on \mathbf{X} , respectively, with $c < n$.
- We denote a $p \times l$ column information matrix on \mathbf{X} by $\mathbf{Z}_\mathbf{X}$ with $l < p$.
- We assume that $\mathbf{Z}_\mathbf{X}$ is also columnwise centered.

External information can assume several forms:

- if the rows of \mathbf{X} represents subjects then we can collect the subject's demographic information in $\mathbf{H}_\mathbf{X}$ and will try to explore how they are related to the set of variables of the main data;
- if we set $\mathbf{H}_\mathbf{X}$ to $\mathbf{1}_n$ then we highlight the mean tendency across the subjects;
- we can consider for $\mathbf{H}_\mathbf{X}$ matrix of dummy variables providing the subject's membership in some prespecified groups in order to analyze the differences among the groups;
- the external information $\mathbf{Z}_\mathbf{X}$ for the variables collected in the matrix \mathbf{X} can be unitary vectors with p component in order to capture the relationships among the columns of the matrix;

- if variables represent stimuli in a preference judgment study then \mathbf{Z}_X can be a matrix of the descriptor variables of the stimuli;
- if we have different within-subjects experimental conditions then \mathbf{Z}_X could be a matrix of contrasts; in this sense, orthogonal contrasts used in Analysis of Variance, by which different linear principal mean effects can be highlighted, are a particular case of external information;
- finally, if variables are repeated observations then \mathbf{Z}_X could be a matrix of coefficients of orthogonal polynomials.

- Let $\mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X} = \mathbf{Z}_X(\mathbf{Z}_X^T \mathbf{Q}_X \mathbf{Z}_X)^- \mathbf{Z}_X^T \mathbf{Q}_X$ be the not symmetric oblique projection operator onto $Im(\mathbf{Z}_X)$ along $Ker(\mathbf{Z}_X^T \mathbf{Q}_X)$ (Takane and Shibayama, 1991).

- "External Analysis" of the Principal Component Analysis with external information on both subjects and variables (Takane and Shibayama, 1991; Takane and Hunter, 2001) considers the additive data decomposition of a single quantitative matrix \mathbf{X} according to the row and column external information matrices \mathbf{H}_X and \mathbf{Z}_X under row and column metrics \mathbf{L} and \mathbf{Q}_X , respectively:

$$\mathbf{X} = \mathbf{P}_{\mathbf{H}_X/\mathbf{L}} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^T + \mathbf{P}_{\mathbf{H}_X/\mathbf{L}}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^T + \mathbf{P}_{\mathbf{H}_X/\mathbf{L}} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^{\perp T} + \mathbf{P}_{\mathbf{H}_X/\mathbf{L}}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^{\perp T} \quad (1)$$

- The four submatrices in (1) are mutually trace-orthogonal in their respective metric matrices \mathbf{L} and \mathbf{Q}_X (Takane and Shibayama, 1991). For example, the first and the second terms in (1) are trace-orthogonal since we have simultaneously

$$\begin{aligned} \underbrace{\text{trace}[(\mathbf{P}_{Z_X/Q_X} \mathbf{X}^T \mathbf{P}_{H_X/L}^T) \mathbf{L}]}_{=\mathbf{A}^T} \underbrace{(\mathbf{P}_{H_X/L}^\perp \mathbf{X} \mathbf{P}_{Z_X/Q_X}^T) \mathbf{Q}_X}_{=\mathbf{B}} &= 0 \\ \underbrace{\text{trace}[\mathbf{L}(\mathbf{P}_{H_X/L} \mathbf{X} \mathbf{P}_{Z_X/Q_X}^T)]}_{=\mathbf{A}} \underbrace{\mathbf{Q}_X(\mathbf{P}_{Z_X/Q_X} \mathbf{X}^T \mathbf{P}_{H_X/L}^\perp)}_{=\mathbf{B}^T} &= 0 \end{aligned}$$

where two matrices of a same size, \mathbf{A} and \mathbf{B} , are said to be trace-orthogonal when $tr(\mathbf{A}^T \mathbf{B}) = tr(\mathbf{A} \mathbf{B}^T) = 0$.

This property implies that the sum of squares of \mathbf{X} can be exactly decomposed according to the sum of sums of squares corresponding to each submatrix in (1):

$$SS(\mathbf{X})_{\mathbf{L}, \mathbf{Q}_X} = SS(\mathbf{P}_{\mathbf{H}_X/\mathbf{L}} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^T)_{\mathbf{L}, \mathbf{Q}_X} + SS(\mathbf{P}_{\mathbf{H}_X/\mathbf{L}}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^T)_{\mathbf{L}, \mathbf{Q}_X} + SS(\mathbf{P}_{\mathbf{H}_X/\mathbf{L}} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^{\perp T})_{\mathbf{L}, \mathbf{Q}_X} + SS(\mathbf{P}_{\mathbf{H}_X/\mathbf{L}}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^{\perp T})_{\mathbf{L}, \mathbf{Q}_X}$$

where $SS(\mathbf{X})_{\mathbf{L}, \mathbf{Q}_X} = tr(\mathbf{X}^T \mathbf{L} \mathbf{X} \mathbf{Q}_X)$.

Each component in (1) has a specific statistical meaning:

- $\mathbf{P}_{\mathbf{H}_X/\mathbf{L}} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^T$ represents the row and column constraints effects;
- $\mathbf{P}_{\mathbf{H}_X/\mathbf{L}}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^T$ reflects the column constraints effect;
- $\mathbf{P}_{\mathbf{H}_X/\mathbf{L}} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^{\perp T}$ reflects the row constraints effect
- the last term $\mathbf{P}_{\mathbf{H}_X/\mathbf{L}}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X/\mathbf{Q}_X}^{\perp T}$ pertains to what can be explained by neither row nor column external information.

Row sides metric matrix \mathbf{L} can assume several forms:

- 1 if there are no differences in importance among the statistical units then \mathbf{L} can be set to the identity matrix;
- 2 if there are differences in importance among the statistical units then special diagonal matrices are used as \mathbf{L} in order to differentially weight the rows of \mathbf{X} ;
- 3 finally, when rows of the data matrix are time points in single-subject multivariate time series data, Escoufier (1987) suggests to use the inverse of the matrix of serial correlations as \mathbf{L} .

- By using this approach, external information \mathbf{H}_X and \mathbf{Z}_X are incorporated within a single data set.
- It is evident that the Takane and Shibayama's basic decomposition (1991) is obtained from (1) by setting $\mathbf{L} = \mathbf{I}_n$ and $\mathbf{Q}_X = \mathbf{I}_p$:

$$\mathbf{X} = \mathbf{P}_{\mathbf{H}_X} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X} + \mathbf{P}_{\mathbf{H}_X}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X} + \mathbf{P}_{\mathbf{H}_X} \mathbf{X} \mathbf{P}_{\mathbf{Z}_X}^\perp + \mathbf{P}_{\mathbf{H}_X}^\perp \mathbf{X} \mathbf{P}_{\mathbf{Z}_X}^\perp$$

- An "Internal analysis" (PCA) is then performed on each component of decomposition.

- Consider the following model

$$\mathbf{X} = \mathbf{H}_X \mathbf{M} \mathbf{Z}_X^T + \mathbf{B} \mathbf{Z}_X^T + \mathbf{H}_X \mathbf{C} + \mathbf{E} \quad (2)$$

where \mathbf{M} ($c \times l$), \mathbf{B} ($n \times l$), and \mathbf{C} ($c \times p$) are matrices of coefficients to be estimated, and \mathbf{E} ($n \times p$) a matrix of error components.

- The four terms in (2) explain portions of the original data matrix, \mathbf{X} .
- The first term pertains to what can be explained by both \mathbf{H}_X and \mathbf{Z}_X , the second term by \mathbf{Z}_X , the third term by \mathbf{H}_X , and the fourth term by neither \mathbf{H}_X nor \mathbf{Z}_X .

- Let $\mathbf{X} = \mathbf{H}_X \mathbf{M} \mathbf{Z}_X^T + \mathbf{E}_1$, and consider the problem of estimating \mathbf{M} so as to minimize $SS(\mathbf{E}_1) = tr(\mathbf{E}_1^T \mathbf{E}_1)$.
- We obtain

$$\hat{\mathbf{M}} = (\mathbf{H}_X^T \mathbf{H}_X)^- \mathbf{H}_X^T \mathbf{X} \mathbf{Z}_X (\mathbf{Z}_X^T \mathbf{Z}_X)^-$$

where $(\mathbf{H}_X^T \mathbf{H}_X)^-$ and $(\mathbf{Z}_X^T \mathbf{Z}_X)^-$ are g-inverses of $(\mathbf{H}_X^T \mathbf{H}_X)$ and $(\mathbf{Z}_X^T \mathbf{Z}_X)$, respectively.

- The residual from the first term is now equal to

$$\hat{\mathbf{E}}_1 = \mathbf{X} - \mathbf{H}_X \hat{\mathbf{M}} \mathbf{Z}_X^T = \mathbf{X} - \mathbf{P}_{H_X} \mathbf{X} \mathbf{P}_{Z_X}$$

where \mathbf{P}_{H_X} and \mathbf{P}_{Z_X} are orthogonal projection operators onto spaces spanned by the column vectors of \mathbf{H}_X and \mathbf{Z}_X , respectively.

- We now separately fit the second and the third terms to $\hat{\mathbf{E}}_1$:

$$\hat{\mathbf{E}}_1 = \mathbf{B}\mathbf{Z}_X^T + \mathbf{E}_2$$

$$\hat{\mathbf{E}}_1 = \mathbf{H}_X\mathbf{C} + \mathbf{E}_3$$

- We obtain a least squares estimate of \mathbf{B} that minimizes $SS(\mathbf{E}_2)$ by

$$\hat{\mathbf{B}} = \mathbf{P}_{\mathbf{H}_X}^\perp \mathbf{X}\mathbf{Z}_X(\mathbf{Z}_X^T\mathbf{Z}_X)^{-}$$

where $\mathbf{P}_{\mathbf{H}_X}^\perp = (\mathbf{I} - \mathbf{P}_{\mathbf{H}_X})$ orthogonal projection operator to $\mathbf{P}_{\mathbf{H}_X}$ (that is $\mathbf{P}_{\mathbf{H}_X}^\perp \mathbf{P}_{\mathbf{H}_X} = \mathbf{0}$ and $\mathbf{P}_{\mathbf{H}_X} + \mathbf{P}_{\mathbf{H}_X}^\perp = \mathbf{I}$).

- Similarly, we obtain

$$\hat{\mathbf{C}} = (\mathbf{H}_X^T\mathbf{H}_X)^{-}\mathbf{H}_X^T\mathbf{X}\mathbf{P}_{\mathbf{Z}_X}^\perp$$

that minimizes $SS(\mathbf{E}_3)$

- Now, the estimate of the fourth term is given by

$$\begin{aligned}\hat{\mathbf{E}} &= \mathbf{H}_X \hat{\mathbf{M}} \mathbf{Z}_X^T + \hat{\mathbf{B}} \mathbf{Z}_X^T + \mathbf{H}_X \hat{\mathbf{C}} \\ &= \mathbf{X} - \mathbf{P}_{H_X} \mathbf{X} \mathbf{P}_{Z_X} - \mathbf{P}_{H_X}^\perp \mathbf{X} \mathbf{P}_{Z_X} - \mathbf{P}_{H_X} \mathbf{X} \mathbf{P}_{Z_X}^\perp \\ &= \mathbf{P}_{H_X}^\perp \mathbf{X} \mathbf{P}_{Z_X}^\perp\end{aligned}$$

- By substituting the least squares estimates for the corresponding parameters, we obtain the following decomposition of the data matrix, \mathbf{X} :

$$\begin{aligned}\mathbf{X} &= (\mathbf{P}_{H_X} + \mathbf{P}_{H_X}^\perp) \mathbf{X} (\mathbf{P}_{Z_X} + \mathbf{P}_{Z_X}^\perp) \\ &= \mathbf{P}_{H_X} \mathbf{X} \mathbf{P}_{Z_X} + \mathbf{P}_{H_X}^\perp \mathbf{X} \mathbf{P}_{Z_X} + \mathbf{P}_{H_X} \mathbf{X} \mathbf{P}_{Z_X}^\perp + \mathbf{P}_{H_X}^\perp \mathbf{X} \mathbf{P}_{Z_X}^\perp\end{aligned}$$