

# Tutorial 5

CHEN Xiao

Department of Mathematics

April 3, 2020

# Principal Component Analysis

Results 5.1 Let  $\Sigma$  be the covariance matrix associated with the random vector  $X' = [X_1, X_2, \dots, X_p]$ . Let  $\Sigma$  have the eigenvalue-eigenvector pair  $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$  where  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ . Then the  $i$ th principal component is given by

$$Y_i = \mathbf{e}_i' \mathbf{X} = e_{i1}X_1 + e_{i2}X_2 + \dots + e_{ip}X_p, i = 1, 2, \dots, p$$

With these choices,

$$\text{Var}(Y_i) = \mathbf{e}_i' \Sigma \mathbf{e}_i = \lambda_i, i = 1, 2, \dots, p$$

$$\text{Cov}(Y_i, Y_k) = \mathbf{e}_i' \Sigma \mathbf{e}_k = 0, i \neq k$$

If some  $\lambda_i$  are equal, the choices of corresponding coefficients vectors,  $\mathbf{e}_i$ , and hence  $Y_i$  are not unique.

Results 5.3 If  $Y_1 = \mathbf{e}'_1 \mathbf{X}$ ,  $Y_2 = \mathbf{e}'_2 \mathbf{X}$ ,  $\dots$ ,  $Y_p = \mathbf{e}'_p \mathbf{X}$  are the principal components obtained from the covariance matrix  $\Sigma$ , then

$$\rho_{Y_i, X_k} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}, i, k = 1, 2, \dots, p$$

are the correlation coefficients between the components  $Y_i$  and the variables  $X_k$ . Here  $(\lambda_1, \mathbf{e}_1), (\lambda_2, \mathbf{e}_2), \dots, (\lambda_p, \mathbf{e}_p)$  are the eigenvalue-eigenvector pair for  $\Sigma$

# Summarizing Sample Variation by Principle Components

Suppose the data  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  represent  $n$  independent drawings from some  $p$ -dimensional population with the mean vector  $\mu$  and covariance matrix  $\Sigma$ . These data yield the sample mean vector  $\bar{\mathbf{x}}$ , the sample covariance matrix  $\mathbf{S}$ , and the sample correlation matrix  $\mathbf{R}$ .

If  $\mathbf{S} = s_{ik}$  be  $p \times p$  sample covariance matrix with eigenvalue-eigenvector pairs  $(\hat{\lambda}_1, \hat{\mathbf{e}}_1), (\hat{\lambda}_2, \hat{\mathbf{e}}_2), \dots, (\hat{\lambda}_p, \hat{\mathbf{e}}_p)$ , the  $i$ th sample principal component is given by  $\hat{y}_i = \hat{\mathbf{e}}_i' \mathbf{x} = \hat{e}_{i1}x_1 + \hat{e}_{i2}x_2 + \dots + \hat{e}_{ip}x_p, i = 1, 2, \dots, p$  where  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$  and  $\mathbf{x}$  is any observation on the variables  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ . Also

Sample variance  $(\hat{y}_k = \hat{\lambda}_k, k = 1, 2, \dots, p)$

Sample covariance  $(\hat{y}_i, \hat{y}_k) = 0, i \neq k$

Total sample variance  $= \sum_{i=1}^n s_{ii} = \hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p$

$$r_{\hat{y}_i, x_k} = \frac{\hat{e}_{ik} \sqrt{\hat{\lambda}_i}}{\sqrt{s_{kk}}}, i, k = 1, 2, \dots, p$$

# The Orthogonal Factor Model

$$X_1 - \mu_1 = \ell_{11}F_1 + \ell_{12}F_2 + \cdots + \ell_{1m}F_m + \varepsilon_1$$

$$X_2 - \mu_2 = \ell_{21}F_1 + \ell_{22}F_2 + \cdots + \ell_{2m}F_m + \varepsilon_2$$

$$\vdots$$

$$X_p - \mu_p = \ell_{p1}F_1 + \ell_{p2}F_2 + \cdots + \ell_{pm}F_m + \varepsilon_p$$

or in matrix notation

$$X - \mu = LF + \varepsilon$$

The coefficient  $\ell_{ij}$  is called the loading of the  $i$ th variable on the  $j$ th factor, so the matrix  $L$  is the matrix of factor loadings.

The unobservable random vectors  $\mathbf{F}$  and  $\varepsilon$  satisfy the following condition  
 $\mathbf{F}$  and  $\varepsilon$  are independent  $E(\mathbf{F}) = 0$ ,  $\text{Cov}(\mathbf{F}) = \mathbf{I}$   $E(\varepsilon) = 0$ ,  $\text{Cov}(\varepsilon) = \Psi$ ,  
where  $\Psi$  is diagonal matrix.

Covariance structure for the Orthogonal Factor Model

1.  $\text{Cov}(\mathbf{X}) = \mathbf{L}\mathbf{L}' + \Psi$

2.  $\text{Cov}(\mathbf{X}, \mathbf{F}) = \mathbf{L}$  or

$$\text{Cov}(X_i, F_j) = \ell_{ij}$$

# The Principal Component Solution of the Factor Model

The principal component analysis of the sample covariance matrix  $S$  is specified in terms of its eigenvalue-eigenvector pairs  $(\hat{\lambda}_1, \hat{\mathbf{e}}_1), (\hat{\lambda}_2, \hat{\mathbf{e}}_2), \dots, (\hat{\lambda}_p, \hat{\mathbf{e}}_p)$ , where  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$ . Let  $m < p$  be the number of common factors. Then the matrix of estimate factor loading  $\{\tilde{\ell}_{ij}\}$  is give by

$$\tilde{\mathbf{L}} = \left[ \sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 : \sqrt{\hat{\lambda}_2} \hat{\mathbf{e}}_2 : \dots : \sqrt{\hat{\lambda}_m} \hat{\mathbf{e}}_m \right]$$

The estimate specific variances are provided by the diagonal elements of the matrix  $\mathbf{S} - \tilde{\mathbf{L}}\tilde{\mathbf{L}}'$ , so

$$\boldsymbol{\Psi} = \begin{bmatrix} \tilde{\psi}_1 & 0 & \cdots & 0 \\ 0 & \tilde{\psi}_2 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & \tilde{\psi}_p \end{bmatrix} \quad \text{with} \quad \tilde{\psi}_i = s_{ii} - \sum_{j=1}^m \tilde{\ell}_{ij}^2$$

## Example 1

(a) Show that the covariance matrix

$$\rho = \begin{bmatrix} 1.00 & 0.63 & 0.45 \\ 0.63 & 1.00 & 0.35 \\ 0.45 & 0.35 & 1.00 \end{bmatrix}$$

for  $p = 3$  standardized random variables  $Z_1, Z_2$  and  $Z_3$  can be generated by the  $m = 1$  factor model  $\mathbf{Z} = \mathbf{LF} + \varepsilon$  with  $\text{Var}(\mathbf{F}) = 1$   $\text{Cov}(\varepsilon, \mathbf{F}) = 0$ , and  $\text{Var}(\varepsilon) = \Psi$

(b) The eigenvalues and eigenvectors of the correlation matrix  $\rho$  above are

$$\begin{aligned} \lambda_1 &= 1.96, & \mathbf{e}'_1 &= [0.625, 0.593, 0.507] \\ \lambda_2 &= 0.68, & \mathbf{e}'_2 &= [-0.219, -0.491, 0.843] \\ \lambda_3 &= 0.36, & \mathbf{e}'_3 &= [0.794, -0.638, -0.177] \end{aligned}$$

Assume an  $m = 1$  factor model, and its loading matrix  $\mathbf{L}$  and matrix of specific variance  $\Psi$  are calculated by the principal component solution method. Calculate  $\text{Corr}(Z_i, F_1)$  for  $i = 1, 2, 3$ .



# Solutions

(a)

$$Z_1 = 0.9F + \varepsilon_1$$

$$Z_2 = 0.7F + \varepsilon_2$$

$$Z_3 = 0.5F + \varepsilon_3$$

Let

$$L = \begin{bmatrix} 0.9 \\ 0.7 \\ 0.5 \end{bmatrix}, \Psi = \begin{bmatrix} 0.19 & 0 & 0 \\ 0 & 0.51 & 0 \\ 0 & 0 & 0.75 \end{bmatrix}$$

Assume  $\text{Var}(F) = 1$ ,  $\text{Cov}(\varepsilon, \mathbf{F}) = 0$ ,  $\text{Var}(\varepsilon) = \Psi$

$$\text{Cov}(\mathbf{Z}) = \rho = \begin{bmatrix} 1.00 & 0.63 & 0.45 \\ 0.63 & 1.00 & 0.35 \\ 0.45 & 0.35 & 1.00 \end{bmatrix}$$

$$\text{Cov}(\mathbf{L}F + \varepsilon) = L \text{Var}(F)L' + 2 \text{Cov}(\varepsilon, \mathbf{F}) + \text{Var}(\varepsilon)$$

$$= LL' + \Psi = \begin{bmatrix} 1.00 & 0.63 & 0.45 \\ 0.63 & 1.00 & 0.35 \\ 0.45 & 0.35 & 1.00 \end{bmatrix} = \rho$$

$$\mathbf{L} = \sqrt{\lambda_1} \mathbf{e}_1 = \begin{bmatrix} 0.8750 \\ 0.8302 \\ 0.7098 \end{bmatrix}$$

The estimate specific variances are provided by the diagonal elements of the matrix  $\boldsymbol{\rho} - \mathbf{L}\mathbf{L}'$ , so

$$\Psi = \begin{bmatrix} 0.234375 & 0 & 0 \\ 0 & 0.310768 & 0 \\ 0 & 0 & 0.496184 \end{bmatrix}$$

$$F_1 = \frac{\mathbf{e}_1' \mathbf{Z}}{\sqrt{\lambda_1}}, \text{Var}(F_1) = \frac{1}{\lambda_1} \mathbf{e}_1' \text{Cov}(\mathbf{Z}) \mathbf{e}_1 = 1$$

$$\text{Cov}(\mathbf{Z}, F_1) = \frac{1}{\sqrt{\lambda_1}} \text{Cov}(\mathbf{Z}, \mathbf{e}_1' \mathbf{Z}) = \frac{1}{\sqrt{\lambda_1}} \mathbf{e}_1' \boldsymbol{\rho} = \sqrt{\lambda_1} \mathbf{e}_1 = \mathbf{L}$$

$$\text{Then } \text{corr}(Z_1, F_1) = \frac{\text{cov}(Z_1, F_1)}{\sqrt{\text{Var}(Z_1) \text{Var} F_1}} = L_1 = 0.8750$$

$$\text{corr}(Z_2, F_1) = \frac{\text{cov}(Z_2, F_1)}{\sqrt{\text{Var}(Z_2) \text{Var} F_1}} = L_2 = 0.8302$$

$$\text{corr}(Z_3, F_1) = \frac{\text{cov}(Z_3, F_1)}{\sqrt{\text{Var}(Z_3) \text{Var} F_1}} = L_2 = 0.7098$$

## Example2

Suppose the random vector  $(X_1, X_2, X_3, X_4, X_5)'$  have the covariance matrix

$$\Sigma = \begin{bmatrix} 1 & 0.75 & 0.75 & 0 & 0 \\ 0.75 & 1 & 0.75 & 0 & 0 \\ 0.75 & 0.75 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & -0.75 \\ 0 & 0 & 0 & -0.75 & 1 \end{bmatrix}$$

(a) Find the first two principal components  $Y_1$  and  $Y_2$ . (b) Find the correlation coefficients  $\rho_{Y_1, X_1}$  and  $\rho_{Y_2, X_4}$ . (c) Find the sum of  $\text{Var}(Y_1) + \text{Var}(Y_2) + \text{Var}(Y_3)$  (d) Determine an appropriate number of principle components.

# Solutions

(a) Let  $\Sigma_1 = \begin{bmatrix} 1 & 0.75 & 0.75 \\ 0.75 & 1 & 0.75 \\ 0.75 & 0.75 & 1 \end{bmatrix}$ ,  $\Sigma_2 = \begin{bmatrix} 1 & -0.75 \\ -0.75 & 1 \end{bmatrix}$  Then

$$\Sigma = \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{bmatrix}$$

$$\det(\Sigma - \lambda I_5) = \det(\Sigma_1 - \lambda I_3) \det(\Sigma_2 - \lambda I_2) = (2.5 - \lambda)(1.75 - \lambda)(0.25 - \lambda)^3$$

Then eigenvalues and eigenvectors are

$$\lambda_1 = 2.50, \lambda_2 = 1.75, \lambda_3 = \lambda_4 = \lambda_5 = 0.25$$

$$\mathbf{e}'_1 = [-0.5773503, -0.5773503, -0.5773503, 0, 0]$$

$$\mathbf{e}'_2 = [0, 0, 0, -0.7071068, 0.7071068]$$

$$\mathbf{e}'_3 = [0.3175691, -0.8102156, 0.4926465, 0, 0]$$

$$\mathbf{e}'_4 = [0.7522078, -0.1010810, -0.6511268, 0, 0]$$

$$\mathbf{e}'_5 = [0, 0, 0, -0.7071068, -0.7071068]$$

The first two principal components

$$Y_1 = \mathbf{e}'_1 X$$

$$Y_2 = \mathbf{e}'_2 X$$

(b)

$$\rho_{Y_1, X_1} = \frac{e_{11}\sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} = \frac{-0.5773503\sqrt{2.5}}{\sqrt{1}} = -0.912871$$

$$\rho_{Y_2, X_4} = \frac{e_{24}\sqrt{\lambda_2}}{\sqrt{\sigma_{44}}} = \frac{-0.7071068\sqrt{1.75}}{\sqrt{1}} = -0.9354144$$

(c)

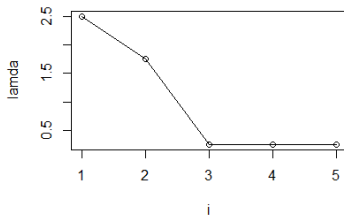
$$\text{Var}(Y_1) + \text{Var}(Y_2) + \text{Var}(Y_3) = \lambda_1 + \lambda_2 + \lambda_3 = 4.5$$

(d)  $\frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4} = 0.5$

$$\frac{\lambda_1 + \lambda_2}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4} = 0.85$$

$$\frac{\lambda_1 + \lambda_2 + \lambda_3}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4} = 0.9$$

A scree plot



## Example3

(a) Show that covariance matrix

$$\rho = \begin{bmatrix} 1.00 & 0.63 & 0.45 & 0.27 \\ 0.63 & 1.00 & 0.35 & 0.21 \\ 0.45 & 0.35 & 1.00 & 0.15 \\ 0.27 & 0.21 & 0.15 & 1.00 \end{bmatrix}$$

for  $p = 4$  standardized random variables  $Z_1, Z_2, Z_3$  and  $Z_4$  can be generated by the  $m = 1$  factor model  $\mathbf{Z} = \mathbf{LF} + \varepsilon$  with  $\text{Var}(\mathbf{F}) = 1$   $\text{Cov}(\varepsilon, \mathbf{F}) = 0$ , and  $\text{Var}(\varepsilon) = \Psi$

(b) The eigenvalues and eigenvectors of the correlation  $\rho$  above are

$$\lambda_1 = 2.0888, \mathbf{e}'_1 = [0.5993, 0.5606, 0.4722, 0.3218]$$

$$\lambda_2 = 0.8832, \mathbf{e}'_2 = [-0.1080, -0.1560, -0.3122, 0.9309]$$

$$\lambda_3 = 0.6739, \mathbf{e}'_3 = [0.2282, 0.5241, -0.8055, -0.1558]$$

$$\lambda_4 = 0.3541, \mathbf{e}'_4 = [-0.7597, 0.6219, 0.1749, 0.0747]$$

Assume that an  $m = 1$  factor model holds, calculate the loading matrix  $L$  and the matrix of the specific variance  $\Psi$  using the principal component solution method. What proportion of the total population variance is explained by the first common factor?

(a)

$$Z_1 = 0.9F + \varepsilon_1$$

$$Z_2 = 0.7F + \varepsilon_2$$

$$Z_3 = 0.5F + \varepsilon_3$$

$$Z_4 = 0.3F + \varepsilon_4$$

Let

$$L = \begin{bmatrix} 0.9 \\ 0.7 \\ 0.5 \\ 0.3 \end{bmatrix}, \Psi = \begin{bmatrix} 0.19 & 0 & 0 & 0 \\ 0 & 0.51 & 0 & 0 \\ 0 & 0 & 0.75 & 0 \\ 0 & 0 & 0 & 0.91 \end{bmatrix}$$

Assume  $\text{Var}(F) = 1$   $\text{Cov}(\epsilon, \mathbf{F}) = 0$ ,  $\text{Var}(\epsilon) = \Psi$

$$\text{Cov}(\mathbf{Z}) = \rho = \begin{bmatrix} 1.00 & 0.63 & 0.45 & 0.27 \\ 0.63 & 1.00 & 0.35 & 0.21 \\ 0.45 & 0.35 & 1.00 & 0.15 \\ 0.27 & 0.21 & 0.15 & 1.00 \end{bmatrix}$$

$$\text{Cov}(\mathbf{L}F + \epsilon) = \mathbf{L} \text{Var}(F) \mathbf{L}' + 2 \text{Cov}(\epsilon, \mathbf{F}) + \text{Var}(\epsilon)$$

$$= \mathbf{L} \mathbf{L}' + \Psi = \begin{bmatrix} 1.00 & 0.63 & 0.45 & 0.27 \\ 0.63 & 1.00 & 0.35 & 0.21 \\ 0.45 & 0.35 & 1.00 & 0.15 \\ 0.27 & 0.21 & 0.15 & 1.00 \end{bmatrix} = \rho$$



# Solutions

$$(b)L = \sqrt{\lambda_1}e_1 = \begin{bmatrix} 0.8661492 \\ 0.8102173 \\ 0.6824556 \\ 0.4650873 \end{bmatrix}$$

The estimate specific variances are provided by the diagonal elements of the matrix  $\rho - LL'$ , so

$$\psi = \begin{bmatrix} 0.2497856 & 0 & 0 & 0 \\ 0 & 0.3435479 & 0 & 0 \\ 0 & 0 & 0.5342543 & 0 \\ 0 & 0 & 0 & 0.7836938 \end{bmatrix}$$

$$\frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4} = 0.5222$$