

Lecture 6: Principal Component Analysis

1 Population Principal Components

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

Consider p random variables X_1, \dots, X_p . Principal component analysis seeks to select a new coordinate system obtained by rotating the original system with X_1, \dots, X_p as the coordinate axes. The new axes represent the directions with maximum variability and provide a simpler and more parsimonious description of the covariance structure.

Let $\mathbf{X}' = (X_1, \dots, X_p)'$. Denote the covariance matrix of \mathbf{X} by Σ whose eigenvalues are $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Let $Y_i = \mathbf{a}'_i \mathbf{X} = a_{i1}X_1 + a_{i2}X_2 + \dots + a_{ip}X_p$ be a linear combination of \mathbf{X} . Then, we have

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

and, more generally, for p such linear combinations,

$$\text{Cov}(Y_i, Y_j) = \mathbf{a}'_i \Sigma \mathbf{a}_j, \quad i, j = 1, \dots, p.$$

The principal components (PC) are those linear combinations Y_1, \dots, Y_p that are uncorrelated and whose variance in (1) are as large as possible. Thus,

The 1st PC = linear combination $\mathbf{a}'_1 \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_1 \mathbf{X})$ subject to $\mathbf{a}'_1 \mathbf{a}_1 = 1$.

The 2nd PC = linear combination $\mathbf{a}'_2 \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_2 \mathbf{X})$ subject to $\mathbf{a}'_2 \mathbf{a}_2 = 1$,

$$\text{and } \text{Cov}(\mathbf{a}'_1 \mathbf{X}, \mathbf{a}'_2 \mathbf{X}) = \mathbf{a}'_2 \Sigma \mathbf{a}_1 = 0.$$

In general,

The i th PC = linear combination $\mathbf{a}'_i \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_i \mathbf{X})$

subject to $\mathbf{a}'_i \mathbf{a}_i = 1$ and $\text{Cov}(\mathbf{a}'_j \mathbf{X}, \mathbf{a}'_i \mathbf{X}) = 0, \quad j = 1, \dots, i - 1$.

Result 8.1. Let Σ be the covariance matrix of the random vector $\mathbf{X}' = (X_1, \dots, X_p)'$. Let the eigenvalue-eigenvector pairs of Σ be $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$ such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$. Then the i th PC of \mathbf{X} is given by

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

Furthermore,

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

If some λ_j are equal then the choices of the corresponding coefficient vectors \mathbf{e}_j and, hence, Y_j are not unique.

Proof. The result follows Eq. (2.51) and (2.52) of the textbook. See Page 80.

Result 8.2. Let $\boldsymbol{\Sigma}$ be the covariance matrix of the random vector $\mathbf{X} = (X_1, \dots, X_p)'$ and $(\lambda_1, \mathbf{e}_1), \dots, (\lambda_p, \mathbf{e}_p)$ be the ordered (in increasing order) eigenvalue-eigenvector pairs of $\boldsymbol{\Sigma}$. Then

$$\sum_{i=1}^p \sigma_{ii} = \sum_{i=1}^p \text{Var}(X_i) = \sum_{i=1}^p \lambda_i = \sum_{i=1}^p \text{Var}(Y_i).$$

The proportion of total population variance explained by the i th PC is

$$\frac{\lambda_i}{\sum_{j=1}^p \lambda_j}.$$

Let $\rho_{X,Y}$ be the correlation coefficient between the random variables X and Y .

Result 8.3. Let $Y_i = \mathbf{e}_i' \mathbf{X}$ be the PC of the random vector \mathbf{X} with covariance matrix $\boldsymbol{\Sigma}$. Then,

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

Proof: Use $X_i = \mathbf{a}_i' \mathbf{X}$, where \mathbf{a}_i is the i th unit vector in R^p , and $\boldsymbol{\Sigma} \mathbf{e}_i = \lambda_i \mathbf{e}_i$.

Correlation matrix: Let $\mathbf{Z} = \mathbf{D}^{-1} \mathbf{X}$, where $\mathbf{D} = \text{diag}\{\sqrt{\sigma_{11}}, \dots, \sqrt{\sigma_{pp}}\}$ with σ_{ii} being the (i,i) th element of the covariance matrix $\boldsymbol{\Sigma}$ of \mathbf{X} . Then the i th PC of \mathbf{Z} is

$$Y_i = \mathbf{e}_i' \mathbf{D}^{-1} (\mathbf{X} - \boldsymbol{\mu}),$$

where $\boldsymbol{\mu} = E(\mathbf{X})$ and $(\lambda_i, \mathbf{e}_i)$ s are the ordered eigenvalue-eigenvector pairs of $\boldsymbol{\rho} = \text{Cov}(\mathbf{Z})$. Moreover,

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

and

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

Remark: Principal component analysis depends on the scales of \mathbf{X} . Thus, in some applications, correlation matrix is used in PCA to obtain scale-invariant results.

2 Sample PC

Suppose that the data $\mathbf{x}_1, \dots, \mathbf{x}_n$ represent n independent draws from some p -dimensional population with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Let $\bar{\mathbf{x}}$, \mathbf{S} and \mathbf{R} be the sample mean, covariance matrix and correlation matrix, respectively. The sample PCs are the counterparts of the population PCs with $\boldsymbol{\Sigma}$ and $\boldsymbol{\rho}$ replaced by \mathbf{S} and \mathbf{R} .

2.1 The number of PCs

The scree plot which is a time-series plot of the eigenvalues of \mathbf{S} or \mathbf{R} in decreasing order. That is, the scatter-plot of $(i, \hat{\lambda}_i)$, $i = 1, \dots, p$. By looking for an elbow (bend) in the scree plot, we can determine the number of PCs. We shall mention some recent works that use information criteria to select the number of principal components later when we discuss factor models.

2.2 Large sample inferences

Asymptotic results for eigenvalues and eigenvectors are only available under the normality assumption and the condition that all eigenvalues are distinct and positive. Some limited results start to appear when the variables are not normally distributed.

Anderson (1963, *Annals of Mathematical Statistics*) derived the following large sample distribution theory for the eigenvalues $\hat{\boldsymbol{\lambda}} = (\hat{\lambda}_1, \dots, \hat{\lambda}_p)'$ and eigenvectors $\hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_p$ of the sample covariance matrix \mathbf{S} :

1. Let $\boldsymbol{\Lambda}$ be the diagonal matrix of eigenvalues $\lambda_1, \dots, \lambda_p$ of $\boldsymbol{\Sigma}$, then $\sqrt{n}(\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}) \sim N_p(\mathbf{0}, 2\boldsymbol{\Lambda}^2)$.

2. Let

$$\mathbf{E}_i = \lambda_i \sum_{k=1, k \neq i}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} \mathbf{e}_k \mathbf{e}_k',$$

then $\sqrt{n}(\hat{\mathbf{e}}_i - \mathbf{e}_i) \sim N_p(\mathbf{0}, \mathbf{E}_i)$.

3. Each $\hat{\lambda}_i$ is distributed independently of the elements of the associated $\hat{\mathbf{e}}_i$.

Property 1 implies that, for large n , the sample eigenvalues $\hat{\lambda}_i$ are independently distributed. Moreover, $\hat{\lambda}_i$ has an approximate $N(\lambda_i, 2\lambda_i^2/n)$ distribution. Using this result, a large sample $100(1 - \alpha)\%$ confidence interval for λ_i is provided by

$$\frac{\hat{\lambda}_i}{1 + z(\alpha/2)\sqrt{2/n}} \leq \lambda_i \leq \frac{\hat{\lambda}_i}{1 - z(\alpha/2)\sqrt{2/n}},$$

where $z(\alpha/2)$ is the upper $100(\alpha/2)$ percentile of a standard normal distribution.

Property 2 implies that the $\hat{\mathbf{e}}_i$'s are normally distributed about the corresponding \mathbf{e}_i 's for large samples. The elements of each $\hat{\mathbf{e}}_i$ are correlated, and the correlation depends to a large

extent on the separation of the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_p$ and the sample size n . Approximate standard errors for the coefficients $\hat{e}_{i,k}$ are given by the square roots of the diagonal elements of $(1/n)\widehat{\mathbf{E}}_i$, where $\widehat{\mathbf{E}}_i$ is derived from \mathbf{E}_i by substituting $\hat{\lambda}_i$'s for λ_i 's and \hat{e}_i 's for the e_i 's.

Testing for the equal correlation structure: Consider the null hypothesis

$$H_o : \boldsymbol{\rho} = \boldsymbol{\rho}_o = \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix}$$

and the alternative hypothesis $H_1 : \boldsymbol{\rho} \neq \boldsymbol{\rho}_o$. Lawley (1963) proposed a testing procedure. Let

$$\bar{r}_k = \frac{1}{p-1} \sum_{i=1, i \neq k}^p r_{ik}, \quad k = 1, 2, \dots, p; \quad \bar{r} = \frac{2}{p(p-1)} \sum_{i < k} r_{ik},$$

$$\hat{\gamma} = \frac{(p-1)^2 [1 - (1 - \bar{r})^2]}{p - (p-2)(1 - \bar{r})^2},$$

where \bar{r}_k is the average of the off-diagonal elements in the k th column of \mathbf{R} and \bar{r} is the overall average of the off-diagonal elements. The large sample approximate α -level test is to reject H_o in favor of \mathbf{H}_a if

$$T = \frac{(n-1)}{(1-\bar{r})^2} \left[\sum_{i < k} \sum (r_{ik} - \bar{r})^2 - \hat{\gamma} \sum_{k=1}^p (\bar{r}_k - \bar{r})^2 \right] > \chi_{(p+1)(p-2)/2}^2(\alpha).$$

Remark: The R command of PC analysis is `princomp`.