

Chap 2. Principal Components Analysis

History

- ▶ First introduced by Karl Pearson (1901) in *Philosophical Magazine* as a procedure for finding lines and planes which best fit a set of points in p -dimensional space. The focus was on **geometric optimization**.

Basic Idea

- The general objectives are
 - ✓ Dimension reduction
 - ✓ Interpretation of data

Reduce the dimensionality of a data set in which there is a large number of inter-related variables while **retaining** as much as possible the variation in the original set of variables.

The reduction is achieved by transforming the original variables to a new set of variables, “**principal components**”, that are uncorrelated and ordered such that the first few retains most of the variation present in the data.

Goals & Objectives

- ▶ Reduction and summary \longrightarrow data reduction.
- ▶ Study the structure of Σ (or \mathbf{S} or \mathbf{R}) \longrightarrow Interpretation.

Study the structure of Σ

illustrates how the overall shape of the data defines the covariance matrix:

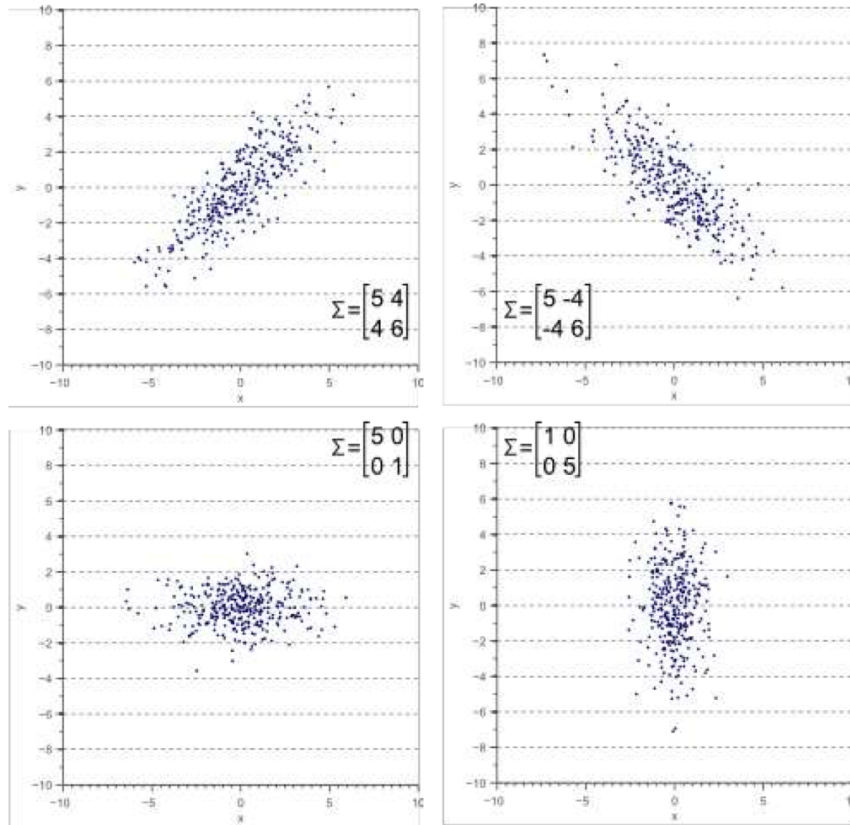


Figure 1. The covariance matrix defines the shape of the data. Diagonal spread is captured by the covariance, while axis-aligned spread is captured by the variance.

Applications

- ▶ Interpretation (study structure)
- ▶ Create a new set of variables (a smaller number that are uncorrelated). These can be used in other procedures (e.g., multiple regression).
- ▶ Select a sub-set of the original variables to be used in other multivariate procedures.
- ▶ Detect outliers or clusters of observations.

- Algebraically, principal components are particular uncorrelated linear combinations of the p random variables X_1, X_2, \dots, X_p .
- Geometrically, principal components represent the selection of a new coordinate system obtained by rotating the original system with X_1, X_2, \dots, X_p as the coordinate axes.

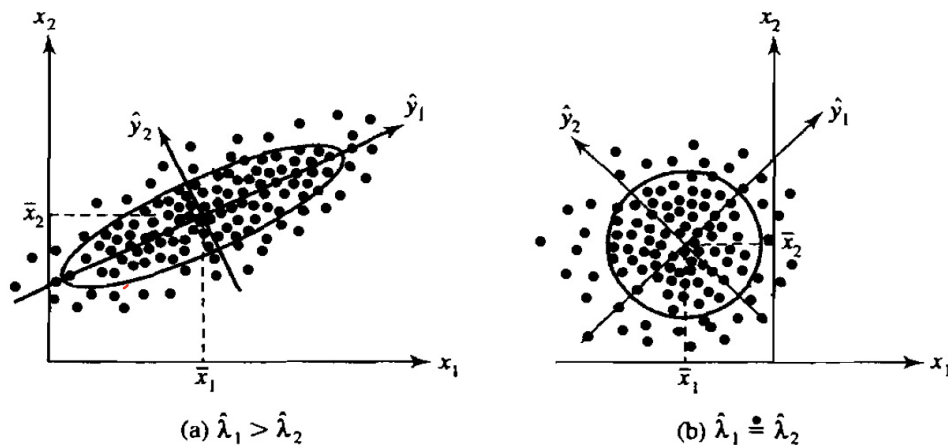


Figure 8.4, Johnson and Wichern (2007)

- ▶ PCs represent a selection of a **new coordinate system** obtained by rotating the original axes to a set of new axes (to provide a simpler structure).
 - ▶ The first principal component represents the direction of maximum variability.
 - ▶ The second principal component represents the direction of maximum variability that is orthogonal to the first.
 - ▶ And so on, until the last PC which represents the direction of minimum variability & orthogonal to all of the others.

4.2 Definition and Derivation of PC's

MATRIX ALGEBRA - REVISIONS

- Principal components represent the directions with maximum variability and provide a simpler and more parsimonious description of the covariance structure.

Let the random vector $\mathbf{X} = (X_1, X_2, \dots, X_p)'$ have the covariance matrix Σ with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$; and $E(\mathbf{X}) = \mathbf{0}$

$$Y_1 = \mathbf{a}'_1 \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$$

$$Y_2 = \mathbf{a}'_2 \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$$

$$\vdots \qquad \qquad \qquad \vdots$$

$$Y_p = \mathbf{a}'_p \mathbf{X} = a_{p1}X_1 + a_{p2}X_2 + \dots + a_{pp}X_p$$

$$\Downarrow$$

$$\text{Var}(Y_i) = \mathbf{a}'_i \Sigma \mathbf{a}_i \qquad \text{Cov}(Y_i, Y_k) = \mathbf{a}'_i \Sigma \mathbf{a}_k$$

Principal components

1st PC = linear combination $\mathbf{a}'_1 \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_1 \mathbf{X})$
 s.t. $\mathbf{a}'_1 \mathbf{a}_1 = 1$

2nd PC = linear combination $\mathbf{a}'_2 \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_2 \mathbf{X})$
 s.t. $\mathbf{a}'_2 \mathbf{a}_2 = 1$ and $\text{Cov}(\mathbf{a}'_1 \mathbf{X}, \mathbf{a}'_2 \mathbf{X}) = 0$

At the i th step,

i th PC = linear combination $\mathbf{a}'_i \mathbf{X}$ that maximizes $\text{Var}(\mathbf{a}'_i \mathbf{X})$
 s.t. $\mathbf{a}'_i \mathbf{a}_i = 1$ and $\text{Cov}(\mathbf{a}'_i \mathbf{X}, \mathbf{a}'_k \mathbf{X}) = 0$, for $k < i$.

Principal components - Result 1:

Let the random vector $\mathbf{X} = (X_1, X_2, \dots, X_p)'$ have the covariance matrix Σ with eigenvalue-eigenvector pairs $(\lambda_1, \underline{\gamma}_1), (\lambda_2, \underline{\gamma}_2), \dots, (\lambda_p, \underline{\gamma}_p)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Then, the i th principal component is

$$Y_i = \underline{\gamma}'_i \mathbf{X} = \gamma_{i1} X_1 + \gamma_{i2} X_2 + \dots + \gamma_{ip} X_p, \quad i = 1, 2, \dots, p$$

and

$$\textcircled{1} \quad \text{Var}(Y_i) = \underline{\gamma}'_i \Sigma \underline{\gamma}_i = \lambda_i \text{ for } i = 1, 2, \dots, p$$

$$\textcircled{2} \quad \text{Cov}(Y_i, Y_k) = \underline{\gamma}'_i \Sigma \underline{\gamma}_k = 0 \text{ for } i \neq k$$

obs: $\underline{\gamma}_i^T \underline{\gamma}_i = 1$ and $\underline{\gamma}_i^T \underline{\gamma}_k = 0$

Norm = 1 and $\underline{\gamma}_i \perp \underline{\gamma}_k$ orthogonal vectors

i -th P.C. : $Y_i = \underline{\gamma}_i^T \mathbf{X}$

observation:

$$Y_i = \underline{\gamma}_i^T (\mathbf{X} - \underline{\mu}) \text{ if}$$

$$E(\mathbf{X}) = \underline{\mu} \neq \underline{0}$$

$$\begin{aligned} \text{Var}(Y_i) &= \\ \text{Var}(\underline{\gamma}_i^T \mathbf{X}) &= \\ \underline{\gamma}_i^T \Sigma \underline{\gamma}_i &= \lambda_i \end{aligned}$$

Prove:

$$\begin{aligned} \frac{\mathbf{a}'\Sigma_X\mathbf{a}}{\mathbf{a}'\mathbf{a}} &= \text{var}(Y_1) \\ \mathbf{a}'\Sigma_X\mathbf{a} &= \text{var}(Y_1)\mathbf{a}'\mathbf{a} \\ \mathbf{a}'\Sigma_X\mathbf{a} - \text{var}(Y_1)\mathbf{a}'\mathbf{a} &= 0 \\ \mathbf{a}'(\Sigma_X\mathbf{a} - \text{var}(Y_1)\mathbf{a}) &= 0 \quad (\text{since } \mathbf{a} \neq 0) \\ \Sigma_X\mathbf{a} - \text{var}(Y_1)\mathbf{a} &= 0 \\ \underbrace{\Sigma_X}_{p \times p} \underbrace{\mathbf{a}}_{p \times 1} &= \underbrace{\text{var}(Y_1)}_{\text{scalar}} \underbrace{\mathbf{a}}_{p \times 1} \end{aligned}$$

which is just the equation what eigenvalues and eigenvectors solve.

So, $Y_i = \gamma_{\tilde{i}}^T X_{\tilde{i}}$ where $\gamma_{\tilde{i}}$ is the i -th eigenvector of Σ .

$$\begin{aligned} \text{Cov}(Y_i, Y_i) &= \text{var}(Y_i) \\ \text{var}(Y_i) &= \text{var}(\gamma_{\tilde{i}}^T X_{\tilde{i}}) = \gamma_{\tilde{i}}^T \text{var}(X_{\tilde{i}}) \gamma_{\tilde{i}} = \\ &= \gamma_{\tilde{i}}^T \underbrace{\Sigma}_{\lambda_i \gamma_{\tilde{i}} \gamma_{\tilde{i}}^T} \gamma_{\tilde{i}} = \lambda_i \gamma_{\tilde{i}}^T \gamma_{\tilde{i}} = \lambda_i \end{aligned}$$

$$\boxed{\Sigma \gamma_{\tilde{i}} = \lambda_i \gamma_{\tilde{i}}}$$

$$\begin{aligned} \text{Cov}(Y_i, Y_k) &= \text{cov}(\gamma_{\tilde{i}}^T X_{\tilde{i}}, \gamma_{\tilde{k}}^T X_{\tilde{k}}) = \gamma_{\tilde{i}}^T \text{cov}(X_{\tilde{i}}, X_{\tilde{k}}) \gamma_{\tilde{k}} = \\ &= \gamma_{\tilde{i}}^T \text{var}(X_{\tilde{i}}) \gamma_{\tilde{k}} = \gamma_{\tilde{i}}^T \underbrace{\Sigma}_{\lambda_k \gamma_{\tilde{k}} \gamma_{\tilde{k}}^T} \gamma_{\tilde{k}} = \\ &= \lambda_k \underbrace{\gamma_{\tilde{i}}^T \gamma_{\tilde{k}}}_{= 0} = 0 \end{aligned}$$

Y_i and Y_k are uncorrelated.

$$\begin{aligned} \gamma_{\tilde{i}}^T \gamma_{\tilde{i}} &= 1 \\ \gamma_{\tilde{i}}^T \gamma_{\tilde{k}} &= 0 \end{aligned}$$

Thus, the PC's are uncorrelated and have variance equal to the eigenvalues of the covariance matrix.

Obs:

1. If some eigenvalues are equal the choice of the corresponding eigenvectors are not unique \Rightarrow PC's is not unique;
2. If $Y_i = \gamma_{\tilde{i}}^T \tilde{X}$ is the i -th P.C., the $Y_i^* = -\gamma_{\tilde{i}}^T \tilde{X}$ is also a P.C. since $(-\gamma_{\tilde{i}}^T)(-\gamma_{\tilde{i}}) = 1$ and $\text{var}(-\gamma_{\tilde{i}}^T \tilde{X}) = \text{var}(\gamma_{\tilde{i}}^T \tilde{X}) = \lambda_i$

Properties of PC + Geometric properties of PC

Principal components - Result 2:

Let the random vector $\mathbf{X} = (X_1, X_2, \dots, X_p)'$ have the covariance matrix Σ with eigenvalue-eigenvector pairs $(\lambda_1, \gamma_1), (\lambda_2, \gamma_2), \dots, (\lambda_p, \gamma_p)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Let $Y_1 = \gamma_1' \mathbf{X}$, $Y_2 = \gamma_2' \mathbf{X}$, \dots , $Y_p = \gamma_p' \mathbf{X}$ be the principal components. Then

$$\sigma_{11} + \dots + \sigma_{pp} = \sum_{i=1}^p \text{Var}(X_i) = \lambda_1 + \dots + \lambda_p = \sum_{i=1}^p \text{Var}(Y_i)$$

Prove:

Matrix Notation:

$$\tilde{Y} = \tilde{T}^T \tilde{X}$$

where

$$\tilde{T} = \begin{matrix} & \gamma_1 & \gamma_2 & \dots & \gamma_p \\ \begin{bmatrix} \gamma_{11} & \gamma_{21} & \dots & \gamma_{p1} \\ \vdots & \vdots & & \vdots \\ \gamma_{1p} & \gamma_{2p} & \dots & \gamma_{pp} \end{bmatrix} \end{matrix}$$

\tilde{T} matrix:

the columns are the eigenvectors of Σ (\tilde{T} is an orthogonal matrix)

$$\text{var}(\underline{Y}) = \text{var}(\underline{T}^T \underline{X}) = \underline{T}^T \text{var}(\underline{X}) \underline{T} = \underline{T}^T \underline{\Sigma} \underline{T}$$

using the spectral decomposition of $\underline{\Sigma}$

$$\underline{\Sigma} = \underline{T} \underline{\Lambda} \underline{T}^T, \text{ we have that}$$

$$\begin{aligned} \text{var}(\underline{Y}) &= \underline{T}^T (\underline{T} \underline{\Lambda} \underline{T}^T) \underline{T} = \underbrace{\underline{T}^T \underline{T}}_{\underline{I}} \underline{\Lambda} \underbrace{\underline{T}^T \underline{T}}_{\underline{I}} = \underline{\Lambda} \\ &= \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix} \end{aligned}$$

$$\text{var}(\underline{X}) = \underline{\Sigma}$$

$$\text{total variance of } \underline{X} = \text{tr}(\underline{\Sigma})$$

$$\text{tr}(\underline{\Sigma}) = \text{tr} \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{12} & \sigma_2^2 & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \sigma_p^2 \end{bmatrix} = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_p^2$$

$$\text{total variance of } \underline{Y} = \text{tr}(\underline{\Lambda})$$

$$\text{tr}(\underline{\Lambda}) = \lambda_1 + \dots + \lambda_p$$

$$\begin{aligned} \text{But: } \text{tr}(\underline{\Sigma}) &= \text{tr}(\underline{T} \underline{\Lambda} \underline{T}^T) = \text{tr}(\underbrace{\underline{T}^T \underline{T}}_{\underline{I}} \underline{\Lambda}) = \\ &= \text{tr}(\underline{\Lambda}) = \lambda_1 + \dots + \lambda_p \end{aligned}$$

$$\text{total variance of } \underline{X} = \text{total variance of } \underline{Y}$$

Generalized variance of $\underline{X} = |\underline{\Sigma}|$

$$|\underline{\Sigma}| = |\underline{P} \underline{\Lambda} \underline{P}^T| = \underbrace{|\underline{P}|}_{\text{orthogonality}} |\underline{\Lambda}| \underbrace{|\underline{P}^T|}_{|\underline{P}| = \pm 1} = |\underline{\Lambda}|$$

$$\text{So, } |\underline{\Sigma}| = |\underline{\Lambda}| = \lambda_1 \times \lambda_2 \times \dots \times \lambda_p$$

thus, the PC's has the same total and generalized variance as the original variables (X 's)

Proportion of total population variance due to k th principal component

$$\frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p}, \quad k = 1, 2, \dots, p$$

Remark: If most (for instance, 80 to 90%) of the total population variance, for large p , can be attributed to the first several principal components, then these components can "replace" the original p variables without much loss of information.

- The magnitude of γ_{ik} measures the importance of the k th variable to the i th principal component, irrespective of the other variables.

Principal components - Result 3:

Let $Y_1 = \gamma_1' \mathbf{X}$, $Y_2 = \gamma_2' \mathbf{X}$, \dots , $Y_p = \gamma_p' \mathbf{X}$ be the principal components obtained from the covariance matrix Σ . Then

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

Prove:

$$\text{Let } \mathbf{e}_k^T = (0, 0, \dots, 1, 0, \dots, 0)$$

k-position

$$X_k = \mathbf{e}_k^T \mathbf{X} \quad ; \quad Y_i = \gamma_i^T \mathbf{X}$$

$$\begin{aligned} \text{Cov}(X_k, Y_i) &= \text{Cov}(\mathbf{e}_k^T \mathbf{X}, \gamma_i^T \mathbf{X}) = \\ &= \mathbf{e}_k^T \text{Cov}(\mathbf{X}, \mathbf{X}) \gamma_i = \mathbf{e}_k^T \underbrace{\Sigma}_{\lambda_i \delta_i} \gamma_i = \\ &= \lambda_i \mathbf{e}_k^T \gamma_i = \lambda_i \delta_{ik} \end{aligned}$$

$$\rho_{Y_i, X_k} = \frac{\text{Cov}(X_k, Y_i)}{\sigma_{X_k} \sigma_{Y_i}} = \frac{\lambda_i \delta_{ik}}{\sigma_k \sqrt{\lambda_i}} = \frac{\gamma_{ik} \sqrt{\lambda_i}}{\sigma_k}$$

$i, k = 1, \dots, p$

- The coefficients γ_{ik} and the correlations ρ_{Y_i, X_k} can lead to different rankings as the measures of the importance of the variables to a given component. However, these rankings are often not appreciably different.
- In practice, variables with relatively large coefficients (in absolute value) tend to have relatively large correlations.
- It is suggested that both the coefficients and the correlations be examined to help interpret the principal components.

Example 1: Let $\underline{x} \in \mathbb{R}^2$ with $\text{Var}(\underline{x}) = \underline{\Sigma} = \begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix}$

obtain the Principal components. $E(\underline{x}) = \underline{\mu}$

$\underline{Y} = \underline{\Gamma}^T \underline{x} \Rightarrow$ obtain the eigenvalues and eigenvectors of $\underline{\Sigma}$

1. Eigenvalues:

$$\text{solution of } |\underline{\Sigma} - \lambda \underline{I}| = 0 \Leftrightarrow \begin{vmatrix} 6 - \lambda & 2 \\ 2 & 3 - \lambda \end{vmatrix} = 0 \Leftrightarrow$$

$$(6 - \lambda)(3 - \lambda) - 4 = 0 \Leftrightarrow 18 - 6\lambda - 3\lambda + \lambda^2 - 4 = 0 \Leftrightarrow$$

$$\lambda^2 - 9\lambda + 14 = 0 \Leftrightarrow \lambda = \frac{9 \pm \sqrt{81 - 54}}{2} \Leftrightarrow \lambda_1 = \frac{9+5}{2} \vee \lambda_2 = \frac{9-5}{2}$$

$$\Leftrightarrow \lambda_1 = 7 \text{ and } \lambda_2 = 2$$

2. Eigenvectors

$$\text{1st: } \underline{\Sigma} \underline{x}_1 = \lambda_1 \underline{x}_1 \Leftrightarrow \begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} = 7 \begin{bmatrix} x_{11} \\ x_{12} \end{bmatrix} \Leftrightarrow$$

$$\begin{cases} 6x_{11} + 2x_{12} = 7x_{11} \\ 2x_{11} + 3x_{12} = 7x_{12} \end{cases} \Leftrightarrow \begin{cases} x_{11} = 2x_{12} \\ x_{12} = 1 \wedge x_{11} = 2 \end{cases} \text{ select:}$$

$$\underline{x}_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad \|\underline{x}_1\| = \sqrt{2^2 + 1^2} = \sqrt{5} \quad \therefore \underline{\gamma}_1 = \begin{bmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix} \quad \frac{\underline{x}_1}{\|\underline{x}_1\|}$$

2nd: $\sum_{\sim} \tilde{x}_{\sim 2} = \lambda_2 \tilde{x}_{\sim 2} \Leftrightarrow$

$$\begin{bmatrix} 6 & 2 \\ 2 & 3 \end{bmatrix} \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} = 2 \begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} \Leftrightarrow \begin{cases} 6x_{21} + 2x_{22} = 2x_{21} \\ 2x_{21} + 3x_{22} = 2x_{22} \end{cases} \Leftrightarrow$$

$$\begin{cases} 4x_{21} = -2x_{22} \Leftrightarrow x_{21} = -\frac{1}{2}x_{22} \end{cases} \text{ select: } x_{22} = 1 \wedge x_{21} = -\frac{1}{2}$$

$$\begin{bmatrix} x_{21} \\ x_{22} \end{bmatrix} = \begin{bmatrix} -1/2 \\ 1 \end{bmatrix} \quad \|\tilde{x}_{\sim 2}\| = \sqrt{(-1/2)^2 + 1^2} = \sqrt{\frac{5}{4}} = \frac{\sqrt{5}}{2}$$

$$\tilde{x}_{\sim 2} = \begin{bmatrix} -\frac{1/2}{\sqrt{5}/2} \\ \frac{1}{\sqrt{5}/2} \end{bmatrix} = \begin{bmatrix} -\frac{1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \end{bmatrix} \quad \tilde{x}_{\sim 2}^T \tilde{x}_{\sim 1} = 0$$

$$Y_1 = \frac{2}{\sqrt{5}}(X_1 - \mu_1) + \frac{1}{\sqrt{5}}(X_2 - \mu_2) \quad \text{var}(Y_1) = \lambda_1 = 7$$

$$Y_2 = -\frac{1}{\sqrt{5}}(X_1 - \mu_1) + \frac{2}{\sqrt{5}}(X_2 - \mu_2) \quad \text{var}(Y_2) = \lambda_2 = 2$$

sum = 9

Obs: $E(Y_1) = E\left[\left(\frac{2}{\sqrt{5}}X_1 - \frac{2}{\sqrt{5}}\mu_1\right) + \left(\frac{1}{\sqrt{5}}X_2 - \frac{1}{\sqrt{5}}\mu_2\right)\right]$
 $= 0$ and $E(Y_2) = 0$

IF $E(\tilde{x}_{\sim}) = \tilde{x}_{\sim}$ then $\begin{cases} Y_1 = \frac{2}{\sqrt{5}}x_1 + \frac{1}{\sqrt{5}}x_2 \text{ with } E(Y_1) = 0 \\ Y_2 = -\frac{1}{\sqrt{5}}x_1 + \frac{2}{\sqrt{5}}x_2 \text{ with } E(Y_2) = 0 \end{cases}$

Example 2: $\underline{X} \in \mathbb{R}^3$ with $E(\underline{X}) = \underline{0}$; $\underline{\Sigma} = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$

with eigenvalues-vectors:

$$\lambda_1 = 5.83; \quad \underline{r}_1^T = [0.383, -0.924, 0]$$

$$\lambda_2 = 2.0; \quad \underline{r}_2^T = [0, 0, 1]$$

$$\lambda_3 = 0.17; \quad \underline{r}_3^T = [0.924, 0.383, 0]$$

$$\underline{X} = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$$

the P.C's are:

$$Y_1 = \underline{r}_1^T \underline{X} = 0.383X_1 - 0.924X_2$$

$$Y_2 = \underline{r}_2^T \underline{X} = X_3$$

$$Y_3 = \underline{r}_3^T \underline{X} = 0.924X_1 + 0.383X_2$$

total variance of \underline{X}

$$= \sigma_1^2 + \sigma_2^2 + \sigma_3^2 =$$

$$+ \text{tr}(\underline{\Sigma}) = 8$$

the random variable X_3 is one of the P.C's because this variable is not correlated with X_1 and with X_2

verify that $\text{var}(\underline{Y}) = \underline{\Lambda} = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix}$

$$\text{var}(aX_1 + bX_2) =$$

$$a^2 \text{var}(X_1) + b^2 \text{var}(X_2) + 2ab \text{cov}(X_1, X_2)$$

$$\text{var}(Y_1) = \text{var}(0.383X_1 - 0.924X_2) = 0.383^2 \text{var}(X_1) + 0.924^2 \text{var}(X_2) - 2 \times 0.383 \times 0.924 \text{cov}(X_1, X_2) =$$

$$= 0.147 \times 1 + 0.854 \times 5 - 0.708(-2) = 5.83 = \lambda_1$$

$$\text{cov}(Y_1, Y_2) = \text{cov}(0.383X_1 - 0.924X_2, X_3) = 0.383 \text{cov}(X_1, X_3) - 0.924 \text{cov}(X_2, X_3) = 0.383 \times 0 - 0.924 \times 0 = 0$$

$$\begin{aligned} \therefore \text{cov}(aX_1 + bX_2, X_3) &= \text{cov}(aX_1, X_3) + \text{cov}(bX_2, X_3) \\ &= a \text{cov}(X_1, X_3) + b \text{cov}(X_2, X_3) \end{aligned}$$

variance of \underline{Y} :

$$\text{var}(\underline{Y}) = \underline{\Lambda} = \begin{bmatrix} 5.83 & 0 & 0 \\ 0 & 2.0 & 0 \\ 0 & 0 & 0.17 \end{bmatrix} \quad \text{tr}(\underline{\Lambda}) =$$

total variance of \tilde{X} is the same of the total variance of \tilde{Y}

$$\text{tr}(\tilde{\Sigma}) = \sigma_1^2 + \sigma_2^2 + \sigma_3^2 = 1 + 5 + 2 = \lambda_1 + \lambda_2 + \lambda_3 = 5.83 + 2 + 0.17 = 8$$

i | % of total variance explained by the i -th P.C.

1 | $\frac{5.83}{8} \times 100\% \approx 72.9\%$

2 | $\frac{2.0}{8} \times 100\% \approx 25\%$

3 | $\frac{0.17}{8} \times 100\% \approx 2.12\%$

} both explained 97.9% of the total variance of \tilde{X}

Correlation between Y_1 and Y_2 with the original variables (X 's)

$$\rho_{Y_1, X_1} = \frac{\gamma_{11} \sqrt{\lambda_1}}{\sigma_1} = \frac{0.383 \sqrt{5.83}}{1} = 0.925$$

$$\rho_{Y_1, X_2} = \frac{\gamma_{12} \sqrt{\lambda_1}}{\sigma_2} = \frac{-0.924 \sqrt{5.83}}{\sqrt{5}} = -0.998$$

$$\rho_{Y_2, X_1} = \rho_{Y_2, X_2} = 0 \text{ and } \rho_{Y_2, X_3} = \frac{\sqrt{2}}{\sqrt{2}} = 1$$

obs: the variable X_2 is the one with greater contribution to the 1st PC and also is the one more correlated with Y_1 .

Principal Component Analysis

Principal components under normality - Geometric Interpretation

Suppose that $\mathbf{X} \sim N_p(\mathbf{0}, \Sigma)$. Then, $\mathbf{x}'\Sigma^{-1}\mathbf{x} = c^2$ is an origin-centered ellipsoid which has axes $\pm c\sqrt{\lambda_i}\mathbf{y}_i$, $i = 1, 2, \dots, p$, where the $(\lambda_i, \mathbf{y}_i)$ are the eigenvalue-eigenvector pairs of Σ . Moreover,

$$\begin{aligned} c^2 = \mathbf{x}'\Sigma^{-1}\mathbf{x} &= \frac{1}{\lambda_1}(\mathbf{y}'_1\mathbf{x})^2 + \frac{1}{\lambda_2}(\mathbf{y}'_2\mathbf{x})^2 + \dots + \frac{1}{\lambda_p}(\mathbf{y}'_p\mathbf{x})^2 \\ &= \frac{1}{\lambda_1}y_1^2 + \frac{1}{\lambda_2}y_2^2 + \dots + \frac{1}{\lambda_p}y_p^2 \end{aligned}$$

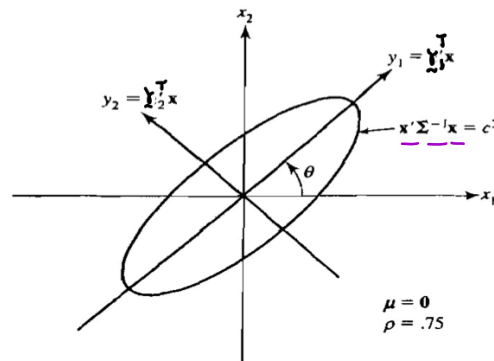


Figure 8.1, Johnson and Wichern (2007)

the PC's lie in the directions of the axes of the constant density ellipsoid.

Principal Components Scores :

objects marks in the P.C. axes

i -th object observations: $(x_{i1}, x_{i2}, \dots, x_{ip})$

i -th object on the j -th P.C.

$$y_{ij} = \gamma_{j1}x_{i1} + \gamma_{j2}x_{i2} + \dots + \gamma_{jp}x_{ip}$$

↓

Score of the i -th object in the j -th P.C

$$\tilde{X} : E(\tilde{X}) = \tilde{0}$$

$$\text{var}(\tilde{X}) = \tilde{\Sigma}$$

Principal Component Analysis

Standardized principal components

Instead of using $\underline{x}^T = (x_1, \dots, x_p)$ and $\text{var}(\underline{x}) = \underline{\Sigma}$

we can calculate p.c. from $\underline{z} = (z_1, \dots, z_p)$, where

$$z_i = \frac{x_i - \mu_i}{\sigma_i} \quad i=1, \dots, p$$

$$\text{var}(\underline{z}) = \underline{\rho}$$

$$\underline{\rho} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \quad \text{tr}(\underline{\rho}) = p$$

Principal Component Analysis

Standardized principal components

$$\underline{\Sigma}^{1/2} = \text{diag}(\sigma_1, \dots, \sigma_p)$$

Let $\mathbf{Z} = (V^{1/2})^{-1}(\mathbf{X} - \boldsymbol{\mu})$. Then $E(\mathbf{Z}) = \mathbf{0}$ and

$$\text{Cov}(\mathbf{Z}) = (V^{1/2})^{-1} \boldsymbol{\Sigma} (V^{1/2})^{-1} = \boldsymbol{\rho}$$

The i th principal component of $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)'$ is given by

$$Y_i = \boldsymbol{\gamma}_i' \mathbf{Z} = \boldsymbol{\gamma}_i' (V^{1/2})^{-1} (\mathbf{X} - \boldsymbol{\mu})$$

and

$$\sum_{i=1}^p \text{Var}(Y_i) = \sum_{i=1}^p \text{Var}(Z_i) = p, \quad \rho_{Y_i, Z_k} = \boldsymbol{\gamma}_{ik} \sqrt{\lambda_i}$$

where the $(\lambda_i, \boldsymbol{\gamma}_i)$ are the eigenvalue-eigenvector pairs of $\boldsymbol{\rho}$

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

The PCs from $\boldsymbol{\Sigma}_X$ are not the same as PCs from $\boldsymbol{\rho}$

We'll look at a situation where standardization makes a difference

This will be the case when the scales of the X variables are (substantially or vastly) different and they are not comparable.

Example 3: Consider the $\underline{X} \in \mathbb{R}^2$; $E(\underline{X}) = \underline{0}$
 covariance Matrix and the Correlation Matrix

$$\underline{\Sigma} = \begin{bmatrix} 1 & 4 \\ 4 & 100 \end{bmatrix}$$

$$\underline{\rho} = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 1 \end{bmatrix}$$

the eigenvalue - eigenvector pairs from

$$\underline{\Sigma} : \begin{cases} \lambda_1 = 100.16; \underline{\gamma}_1^T = [0.040; 0.999] \\ \lambda_2 = 0.84; \underline{\gamma}_2^T = [0.999; -0.040] \end{cases}$$

$$\underline{\rho} : \begin{cases} \lambda_1^* = 1.4; \underline{\gamma}_{\rho 1}^{*T} = [0.707; 0.707] \\ \lambda_2^* = 0.6; \underline{\gamma}_{\rho 2}^{*T} = [0.707; -0.707] \end{cases}$$

$$z_1 = \frac{X_1 - 0}{\sigma_1}$$

the P.C.'s are using

$$\underline{\Sigma} : \begin{cases} Y_1 = 0.040 X_1 + 0.999 X_2 \\ Y_2 = 0.999 X_1 - 0.040 X_2 \end{cases} \quad \underline{\rho} : \begin{cases} Y_1^* = 0.707 z_1 + 0.707 z_2 \\ Y_2^* = 0.707 z_1 - 0.707 z_2 \end{cases}$$

• We see that X_2 completely dominates the 1st P.C. of $\underline{\Sigma}$
 this component Y_1 explains $\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{100.16}{101} = 0.992$ of the total
 variance of \underline{X}

• In contrast, the variables z_1 and z_2 contribute equally to the
 P.C.'s of $\underline{\rho}$

In this case the 1st P.C. explains $\frac{\lambda_1^*}{\rho} = \frac{1.4}{2} = 0.7$ of the
 total variance of the standardized variables (\underline{z})

$$\rho_{Y_1, X_1} = \frac{\sigma_{11} \sqrt{\lambda_1}}{\sigma_1} = \frac{0.04 \sqrt{100.16}}{1} = 0.4$$

$$\rho_{Y_1, X_2} = \frac{\sigma_{12} \sqrt{\lambda_1}}{\sigma_2} = \frac{0.999 \sqrt{100.16}}{10} \approx 0.999$$

$$\rho_{Y_1^*, z_1} = \sigma_{11}^* \sqrt{\lambda_1^*} = 0.707 \sqrt{1.4} = 0.837$$

$$\rho_{Y_1^*, z_2} = \sigma_{12}^* \sqrt{\lambda_1^*} = 0.837$$

We conclude that the relative importance of the variables is affected by standardization.

conclusion: the P.C.'s of $\tilde{\mathbf{z}}$ differ from those of \mathbf{z} if the variables are often standardized when they have different units or widely different scales.

Sample principal components

$$\hat{\Sigma} = \tilde{\Sigma} \quad \hat{\mu} = \tilde{\mu}$$

Used to summarize the sample variation by PCs.

The Algebra is the same as in population principal components.

- ▶ $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are n independent observations from a population with μ and Σ .
- ▶ $\bar{\mathbf{x}}_{p \times 1}$ = sample mean vector.
- ▶ $\mathbf{S}_{p \times p} = \{s_{ik}\}$ = sample covariance matrix.
- ▶ \mathbf{S} has eigenvalue/vector pairs $(\hat{\lambda}_1, \hat{\mathbf{y}}_1), \dots, (\hat{\lambda}_p, \hat{\mathbf{y}}_p)$ where $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$.
- ▶ The $\hat{\cdot}$ indicates these are estimates of population values.
- ▶ The i^{th} sample principal component is given by

$$\hat{\mathbf{y}}_i = \hat{\mathbf{y}}_i^T \mathbf{x} = \hat{y}_{i1}x_1 + \hat{y}_{i2}x_2 + \dots + \hat{y}_{ip}x_p$$

- ▶ The i^{th} PC sample variance = $\widehat{\text{var}}(\hat{y}_i) = \hat{\lambda}_i$ for $i = 1, \dots, p$.
- ▶ The PC sample covariances = $\widehat{\text{cov}}(\hat{y}_i, \hat{y}_k) = 0$ for all $i \neq k$.

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

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

\uparrow
 $\hat{\rho}_{\hat{y}_i, x_k} = r_{\hat{y}_i, x_k}$

Standardized sample principal components

Let $\mathbf{z}_j = D^{-1/2}(\mathbf{x}_j - \bar{\mathbf{x}})$, $j = 1, 2, \dots, p$, be the standardized observations with covariance matrix $R = (D^{1/2})^{-1}S(D^{1/2})^{-1}$. Then, the i th sample principal component is given by

$$\hat{y}_i = \mathbf{v}_i' \mathbf{z} = \hat{\mathbf{v}}_{i1}z_1 + \hat{\mathbf{v}}_{i2}z_2 + \dots + \hat{\mathbf{v}}_{ip}z_p, \quad i = 1, 2, \dots, p$$

where the $(\hat{\lambda}_i, \hat{\mathbf{v}}_i)$ are the eigenvalue-eigenvector pairs of R with $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$. Moreover,

$$\text{Sample variance}(\hat{y}_i) = \hat{\lambda}_i, \quad i = 1, 2, \dots, p$$

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

$$\text{Total standardized sample variance} = \text{tr}(R) = p = \sum_{i=1}^p \hat{\lambda}_i$$

$$r_{\hat{y}_i, z_k} = \hat{\mathbf{v}}_{ik} \sqrt{\hat{\lambda}_i}$$

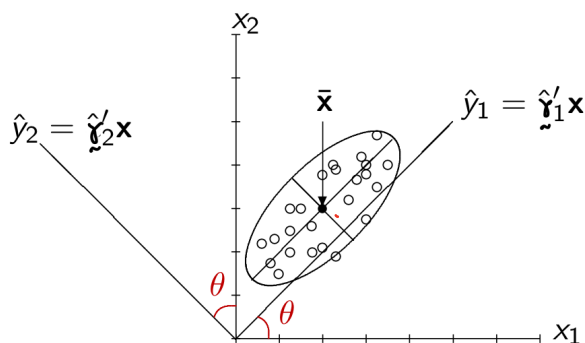
- the sample P.C. based on $\hat{\Sigma}$ are not the same of those based on \hat{R}
- the variable are often standardized when they have different units and variances.

Interpretation of sample principal components

- ▶ PCs based on a sample of n p -dimensional observations are new variables specified by a rigid rotation of the original axes to a new orientation such that the directions of the axes in the new orientation have maximum variances in the sample.
 - ▶ The rotation must be **rigid** since the new variables must be \perp .
 - ▶ Directions of the new axes are based on **S** (or **R**)

The centered sample principal components $\hat{y}_i = \hat{\gamma}'_i(\mathbf{x} - \bar{\mathbf{x}})$, $i = 1, 2, \dots, p$, can be viewed as the result of translating the origin of the original coordinate system to $\bar{\mathbf{x}}$ and then rotating the coordinate axes until they pass through the scatter in the directions of maximum variance.

Geometry of Sample PC



The PCs are projections of observations onto the principal axes of the ellipsoids.

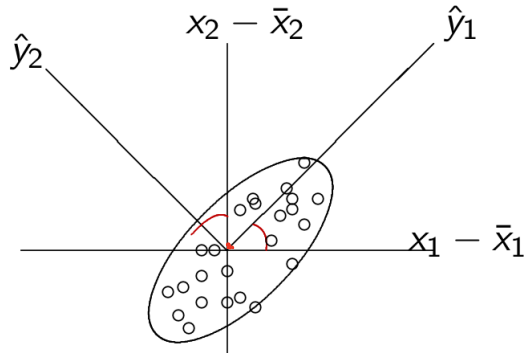
We can re-center the x 's, which also centers the \hat{y} 's; that is

$$(\mathbf{x}_i - \bar{\mathbf{x}}) = 0 \longrightarrow \hat{y}_i \text{ has mean } 0$$

Subtraction of $\bar{\mathbf{x}}$ only effects the mean and does not effect variances and covariance.

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \xrightarrow{\text{shift location}} \begin{pmatrix} x_1 - \bar{x}_1 \\ x_2 - \bar{x}_2 \end{pmatrix} \xrightarrow{\text{rigid rotation}} \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix}$$

The PCs are projections of observations onto the principal axes of the ellipsoids.



How Many Components to Retain ??

The number of principal components

How many principal components should be retained? (No definite answer)

- The amount of total sample variance explained;
- The relative sizes of the eigenvalues;
- The subject-matter interpretations of the components

✓ A useful visual aid to determining an appropriate number of principal components is a *scree plot* (the magnitude of an eigenvalue vs. its number).

Remark: A component associated with an eigenvalue near zero and, hence, deemed unimportant, may indicate an unsuspected linear dependency in the data.

A scree plot

The number of components is taken to be the point at which the remaining eigenvalues are relatively small and all about the same size.

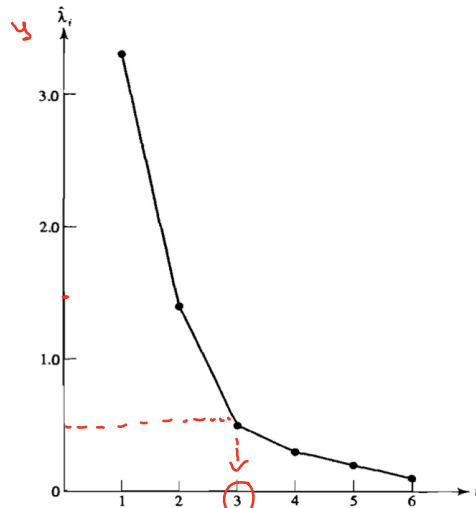


Figure 8.2, Johnson and Wichern (2007)

Some possible rules for choosing the number of P.C.'s :

- retain the first k P.C.'s to explain $\approx 80\% - 90\%$ of the total variance;

- retain P.C.'s with eigenvalues $\geq \bar{\lambda} = \frac{\sum_{i=1}^p \hat{\lambda}_i}{p}$

using the standardized variables $\bar{\lambda} = 1$

Kaiser rule : keep components with $\hat{\lambda} \geq 1$

Graphing Principal Components

- Plots of the principal components can reveal suspect observations, as well as provide checks on the assumption of normality.

- ▶ Reveal **suspect** observations (outliers, influential observations).
- ▶ Check multivariate normality assumptions.
- ▶ Look for clusters.
- ▶ Provide insight into structure in the data.

Suspect Observations

- ▶ The first PCs can help reveal **influential** observations: those that contribute more to variances than other observations such that if we removed them the results change quite a bit.
- ▶ The last PCs can help to reveal **outliers**: those observations that are a typical of the data set; they're inconsistent with the rest of the data (could be miss-coded).

Graphing the principal components

- To help check the normal assumption, construct scatter diagrams for pairs of the first few principal components. Also, make Q-Q plots from the sample values generated by *each* principal components.
- Construct scatter diagrams and Q-Q plots for the last few principal components. These help identify suspect observations.

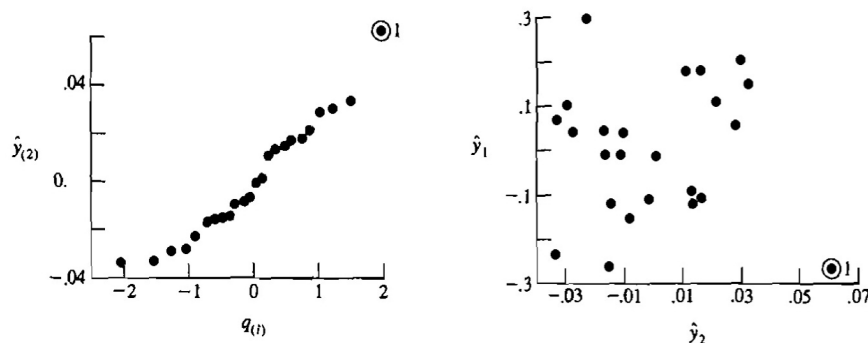


Figure 8.5 and 8.6, Johnson and Wichern (2007)

P.e's Interpretation (not easy task)

Example 4. The weekly rates of return for five stocks listed on the New York Stock Exchange were determined for the period January 1975 through December 1976. Let x_1, x_2, \dots, x_5 denote observed weekly rates of return for the five stocks. Then

$$\bar{x} = [0.0054, 0.0048, 0.0057, 0.0063, 0.0037]^T$$

and

$$R = \begin{bmatrix} 1.000 & 0.577 & 0.509 & 0.387 & 0.462 \\ 0.577 & 1.000 & 0.599 & 0.389 & 0.322 \\ 0.509 & 0.599 & 1.000 & 0.436 & 0.426 \\ 0.387 & 0.389 & 0.436 & 1.000 & 0.523 \\ 0.462 & 0.322 & 0.426 & 0.523 & 1.000 \end{bmatrix}$$

Different
Notation:

The eigenvalues and corresponding normalized eigenvectors of R are:

$$\hat{\lambda}_1 = 2.857 \quad \hat{e}_1 = [0.464, 0.457, 0.470, 0.421, 0.421]^T$$

$$\hat{\lambda}_2 = 0.809 \quad \hat{e}_2 = [0.240, 0.509, 0.260, -0.526, -0.582]^T$$

$$\hat{\lambda}_3 = 0.540 \quad \hat{e}_3 = [-0.612, 0.178, 0.335, 0.541, -0.435]^T$$

$$\hat{\lambda}_4 = 0.452 \quad \hat{e}_4 = [0.387, 0.206, -0.662, 0.472, -0.382]^T$$

$$\hat{\lambda}_5 = 0.343 \quad \hat{e}_5 = [-0.451, 0.676, -0.400, -0.176, 0.385]^T$$

$$\hat{e}_i = \hat{v}_i$$

Using the standardized variables, we obtain the first two sample principal components

$$\hat{y}_1 = \hat{e}_1^T \mathbf{z} = 0.464z_1 + 0.457z_2 + 0.470z_3 + 0.421z_4 + 0.421z_5$$

$$\hat{y}_2 = \hat{e}_2^T \mathbf{z} = 0.240z_1 + 0.509z_2 + 0.260z_3 - 0.526z_4 - 0.582z_5$$

These components account for

$$\left(\frac{\hat{\lambda}_1 + \hat{\lambda}_2}{p} \right) \times 100\% = 73\%$$

of the total sample variance. The first component is an equally weighted sum, or “index”, of the five stocks. This component might be called a market component. The second component represents a contrast between the first three stocks (which were chemical stocks) and the last two stocks (oil stocks). It might be called an industry component.

Example 5: % of People employed in different industries

Dataset: $\begin{cases} 26 \text{ countries} \\ 3 \text{ variables} \end{cases}$

X_1 = 'percent in manufacturing'

X_2 = ' " in services industry'

X_3 = ' " in social and personal services'

$$\mathbf{S} = \begin{bmatrix} 49.109 & 6.535 & 7.379 \\ & 20.933 & 17.876 \\ & & 46.643 \end{bmatrix}, \text{ with eigenvectors} \\ \text{and eigenvalues:}$$

$$\hat{\gamma}_1^T = (0.580; 0.396; 0.712) ; \hat{\lambda}_1 = 62.62$$

$$\hat{\gamma}_2^T = (0.811; -0.207; -0.546) ; \hat{\lambda}_2 = 42.47$$

$$\hat{\gamma}_3^T = (-0.069; 0.894; -0.442) ; \hat{\lambda}_3 = 11.60$$

P.C.S:

$$\hat{Y}_1 = 0.580(X_1 - \bar{x}_1) + 0.396(X_2 - \bar{x}_2) + 0.712(X_3 - \bar{x}_3)$$

$$\hat{Y}_2 = 0.811(X_1 - \bar{x}_1) - 0.207(X_2 - \bar{x}_2) - 0.546(X_3 - \bar{x}_3)$$

$$\hat{Y}_3 = -0.069(X_1 - \bar{x}_1) + 0.894(X_2 - \bar{x}_2) - 0.442(X_3 - \bar{x}_3)$$

$$\text{var}(\hat{Y}_1) = \hat{\lambda}_1$$

variance explained:

i	$\widehat{\text{Var}}(\hat{Y}_i) = \hat{\lambda}_i$	% cumulative variance explained = $\frac{\sum_{j=1}^i \lambda_j}{\sum_{j=1}^3 \lambda_j} \times 100\%$
1	62.62	$\frac{\hat{\lambda}_1}{\hat{\lambda}_1 + \hat{\lambda}_2 + \hat{\lambda}_3} \times 100\% = \frac{62.62}{116.68} \times 100\% = 53.66$
2	42.47	90.06
3	11.60	100.00
total	116.68	

$$\text{tr}(\hat{\Sigma}) = 49.109 + 20.933 + 46.643$$

Correlations: $r_{\hat{Y}_i, X_k} = \frac{\hat{\gamma}_{ik} \sqrt{\hat{\lambda}_i}}{\lambda_k}, i, k = 1, \dots, 3$

original variables	\hat{Y}_1	\hat{Y}_2	\hat{Y}_3
X_1	$\frac{\hat{\gamma}_{11} \sqrt{\hat{\lambda}_1}}{\lambda_1} = \frac{\sqrt{62.62 \times 0.580}}{\sqrt{49.109}} = 0.66$	0.75	
X_2	$\frac{0.396 \sqrt{62.62}}{\sqrt{20.933}} = 0.69$	-0.30	
X_3	0.82	-0.52	

Possible Interpretation of P.C's:

- \hat{Y}_1 : all variables are contributing to the 1st P.C., it could be an overall percent employment in all industries (global measure)
- \hat{Y}_2 : Is a contrast between manufacturing (X_1) with service (X_2) and social (X_3)
- \hat{Y}_3 : Contrast between Industry (X_2) and Social (X_3)

P.C's Explanation:

- Based on Loadings ($\tilde{\gamma}_i \rightarrow$ weights of variables X 's)
- Based on correlations between PC's and the variables (X 's or Z 's)