

Chapter 8 - Principal Components

8.1

Goal: Describe variance-covariance structure of set of variables through linear combinations of the variables. 1) Data Reduction, 2) Interpretation.

n measurements on p variables \Rightarrow n measurements on $k \leq p$ components.

8.2 Population Principal Components

p random variables: X_1, \dots, X_p

or: $X = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix}$ w/ covariance matrix Σ
w/ eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$

Linear combinations:

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

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

\vdots

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

$$V\{Y_i\} = \underline{a}_i' V\{X\} \underline{a}_i = \underline{a}_i' \Sigma \underline{a}_i \quad i=1, \dots, p$$

$$\text{Cov}\{Y_i, Y_k\} = \underline{a}_i' \Sigma \underline{a}_k \quad i, k=1, \dots, p$$

- First principal component is the linear combination $\underline{a}_1' X$ that maximizes $V\{a_1' X\}$ s.t. $\underline{a}_1' \underline{a}_1 = 1$
- Second principal component $\equiv \underline{a}_2' X$ that maximizes $V\{a_2' X\}$ s.t. $\underline{a}_2' \underline{a}_2 = 1$ $\text{Cov}\{\underline{a}_1' X, \underline{a}_2' X\} = 0$. . .

Result 8.1 Random vector w/ $\underline{\tilde{X}} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix}$ and $V\{\underline{\tilde{X}}\} = \underline{\tilde{\Sigma}}$

$\underline{\tilde{\Sigma}}$ has eigenvalue - eigenvector ~~pairs~~ pairs:

$$(\lambda_1, \underline{e}_1), (\lambda_2, \underline{e}_2), \dots, (\lambda_p, \underline{e}_p) \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

Then the i^{th} principal component is given as:

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

$$\text{with } V\{Y_i\} = \underline{e}_i' \underline{\tilde{\Sigma}} \underline{e}_i = \lambda_i \quad i=1, \dots, p$$

$$\text{Cov}\{Y_i, Y_k\} = \underline{e}_i' \underline{\tilde{\Sigma}} \underline{e}_k = 0 \quad i \neq k$$

Proof: Equation (2-51): $B = \underline{\tilde{\Sigma}}$

$$\max_{\underline{a} \neq 0} \frac{\underline{a}' \underline{\tilde{\Sigma}} \underline{a}}{\underline{a}' \underline{a}} = \lambda_1 \quad \text{attained @ } \underline{a} = \underline{e}_1 \quad (\underline{e}_1' \underline{e}_1 = 1) \quad \text{Normalized} \Rightarrow$$

$$\Rightarrow \max_{\underline{a} \neq 0} \frac{\underline{a}' \underline{\tilde{\Sigma}} \underline{a}}{\underline{a}' \underline{a}} = \lambda_1 = \frac{\underline{e}_1' \underline{\tilde{\Sigma}} \underline{e}_1}{\underline{e}_1' \underline{e}_1} = \underline{e}_1' \underline{\tilde{\Sigma}} \underline{e}_1 = V\{Y_1\}$$

Equation (2-52):

$$\max_{\underline{a} \perp \underline{e}_1, \dots, \underline{e}_k} \frac{\underline{a}' \underline{\tilde{\Sigma}} \underline{a}}{\underline{a}' \underline{a}} = \lambda_{k+1} \quad @ \quad \underline{a} = \underline{e}_{k+1} \quad k=1, \dots, p-1$$

$$\underline{e}_{k+1}' \underline{e}_i = 0 \quad i=1, \dots, k, \quad k=1, \dots, p-1$$

$$\frac{\underline{e}_{k+1}' \underline{\tilde{\Sigma}} \underline{e}_{k+1}}{\underline{e}_{k+1}' \underline{e}_{k+1}} = \underline{e}_{k+1}' \underline{\tilde{\Sigma}} \underline{e}_{k+1} = V\{Y_{k+1}\}$$

$$\text{where: } \underline{e}_{k+1}' (\underline{\tilde{\Sigma}} \underline{e}_{k+1}) = \lambda_{k+1} \underline{e}_{k+1}' \underline{e}_{k+1} = \lambda_{k+1} (1) = V\{Y_{k+1}\}$$

- If all eigenvalues are distinct \Rightarrow all eigenvectors orthogonal
- If eigenvalues are not all distinct, eigenvectors corresponding to common eigenvalues can be chosen to be orthogonal.

$$\Rightarrow \underline{e}_i' \underline{e}_k = 0 \quad i \neq k$$

$$\text{Cov}\{Y_i, Y_k\} = \underline{e}_i' \underline{\Sigma} \underline{e}_k = \underline{e}_i' \lambda_k \underline{e}_k = \lambda_k \underline{e}_i' \underline{e}_k = 0 \quad i \neq k$$

Result 8.2 R.V. $\underline{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix}$ w/ $v\{X\} = \underline{\Sigma}$ w/

eigenvalue/eigenvector pairs $(\lambda_1, \underline{e}_1), (\lambda_2, \underline{e}_2), \dots, (\lambda_p, \underline{e}_p)$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

$Y_i = \underline{e}_i' \underline{X}, \dots, Y_p = \underline{e}_p' \underline{X}$ = principal components

Then: $\sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \sum_{i=1}^p v\{X_i\} = \lambda_1 + \dots + \lambda_p = \sum_{i=1}^p v\{Y_i\}$

Proof: $\sigma_{11} + \dots + \sigma_{pp} = \text{tr}(\underline{\Sigma})$ w/ $\underline{\Sigma} = P \Lambda P'$

$$P = [\underline{e}_1 \dots \underline{e}_p] \quad \Lambda = \text{diag}\{\lambda_i\} \quad PP' = P'P = I$$

$$\text{tr}(\underline{\Sigma}) = \text{tr}(P \Lambda P') = \text{tr}(P' P \Lambda) = \text{tr}(\Lambda) = \sum_{i=1}^p v\{Y_i\}$$

Proportion of total population variance due to k^{th} Prin. Comp.:

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

If most of variation for large p can be explained by a few principal components, data can be reduced w/out ^{major} loss of information.

Result 8.3 $Y_1 = \underline{e}_1' X, \dots, Y_p = \underline{e}_p' X$ PC's from Σ

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

Proof: Let $\underline{a}_k' = [0 \dots 0 \ 1 \ 0 \dots 0] \Rightarrow X_k = \underline{a}_k' X$

$\text{Cov}\{\underline{a}_k' X, \underline{e}_i' X\} = \underline{a}_k' \Sigma \underline{e}_i \quad \Sigma \underline{e}_i = \lambda_i \underline{e}_i$

$\Rightarrow \text{Cov}\{X_k, Y_i\} = \underline{a}_k' \lambda_i \underline{e}_i = \lambda_i e_{ik}$

$V\{Y_i\} = \lambda_i \quad V\{X_k\} = \sigma_{kk} \Rightarrow \text{Corr}\{X_k, Y_i\} = \frac{\lambda_i e_{ik}}{\sqrt{\lambda_i} \sqrt{\sigma_{kk}}} = \frac{e_{ik} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}}$

EXAMPLE - LPGA DATA ($N=156, P=4$ (Drive, Fault, Greens, Putts))

$\Sigma = \begin{bmatrix} 86.23 & -25.41 & 15.60 & -1.72 \\ & 31.16 & 4.60 & 1.12 \\ & & 12.95 & 0.20 \\ & & & 1.06 \end{bmatrix}$

$\text{trace}(\Sigma) = 131.40 \quad \lambda_1 = 98.18 \quad \lambda_3 = 5.02 \quad \sum_{i=1}^4 \lambda_i = 131.41$
 $\lambda_2 = 27.21 \quad \lambda_4 = 1.00$

$$\tilde{e}_1 = \begin{pmatrix} 0.927 \\ -0.342 \\ 0.151 \\ -0.020 \end{pmatrix}$$

$$\tilde{e}_2 = \begin{pmatrix} -0.225 \\ -0.829 \\ -0.512 \\ -0.025 \end{pmatrix}$$

$$\tilde{e}_3 = \begin{pmatrix} -0.299 \\ -0.442 \\ 0.844 \\ 0.048 \end{pmatrix}$$

$$\tilde{e}_4 = \begin{pmatrix} 0.027 \\ -0.007 \\ -0.049 \\ 0.998 \end{pmatrix}$$

- Component 1: Driving Distance (Note that σ_{11} is large)
- Component 2: Fairway Put/Greens in Reg
- Component 3: Greens in Reg.
- Component 4: Putting.

$$\frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4} = \frac{98.18}{131.41} = .747$$

$$\frac{\lambda_1 + \lambda_2}{131.41} = \frac{98.18 + 27.21}{131.41} = .954$$

First 2 components could replace 4 vars w/
little loss of information.

$$R_{Y_1, X_1} = \frac{e_{11} \sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} = \frac{0.927 \sqrt{98.18}}{\sqrt{86.23}} = .989$$

$$R_{Y_1, X_2} = \frac{e_{12} \sqrt{\lambda_1}}{\sqrt{\sigma_{22}}} = \frac{-0.342 \sqrt{98.18}}{\sqrt{31.16}} = -.607$$

Case where $\underline{X} \sim N_p(\underline{\mu}, \underline{\Phi})$

Density of X constant on $\underline{\mu}$ centred ellipsoids:

$$(\underline{x} - \underline{\mu})' \underline{\Phi}^{-1} (\underline{x} - \underline{\mu}) = c^2 \quad \text{w/ axes } \pm c\sqrt{\lambda_i} \underline{e}_i \\ i=1, \dots, p$$

Setting $\underline{\mu} = 0$ ($\underline{W} = \underline{X} - \underline{\mu}$) $\Rightarrow E\{W\} = 0$, $\text{Cov}\{W\} = \text{Cov}\{X\}$

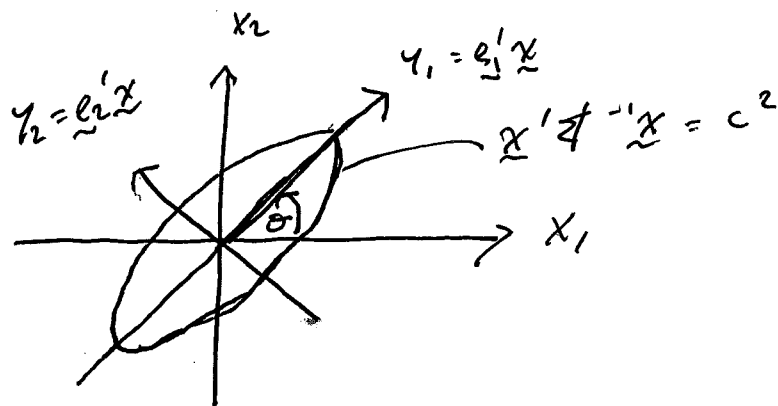
$$c^2 = \underline{x}' \underline{\Phi}^{-1} \underline{x} = \frac{1}{\lambda_1} (\underline{e}_1' \underline{x})^2 + \dots + \frac{1}{\lambda_p} (\underline{e}_p' \underline{x})^2$$

$$\text{Setting } y_i = \underline{e}_i' \underline{x} \Rightarrow c^2 = \frac{1}{\lambda_1} y_1^2 + \dots + \frac{1}{\lambda_p} y_p^2$$

($\lambda_1, \dots, \lambda_p > 0$)

$y_1 = \underline{e}_1' \underline{x}$, ..., $y_p = \underline{e}_p' \underline{x}$ lie in directions of axes of a

constant density ellipsoid.



Principal Components from Standardized Variables 8.7

$$z_1 = \frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}}, \dots, z_p = \frac{x_p - \mu_p}{\sqrt{\sigma_{pp}}}$$

$$\underline{z} = (V^{1/2})^{-1} (\underline{x} - \underline{\mu}) \quad V^{1/2} = \begin{bmatrix} \sqrt{\sigma_{11}} & & 0 \\ & \sqrt{\sigma_{22}} & \\ 0 & & \dots \\ & & & \sqrt{\sigma_{pp}} \end{bmatrix}$$

$$E\{\underline{z}\} = (V^{1/2})^{-1} [\underline{\mu} - \underline{\mu}] = \underline{0}$$

$$V\{\underline{z}\} = (V^{1/2})^{-1} V (V^{1/2})^{-1} = \underline{I}$$

Result 8.4 i^{th} principal component of standardized variable $\underline{z}' = (z_1 \dots z_p)$ w/ $\text{cov}\{\underline{z}\} = \underline{I}$ is:

$$Y_i = \underline{e}_i' \underline{z} = \underline{e}_i' (V^{1/2})^{-1} (\underline{x} - \underline{\mu}) \quad i=1, \dots, p$$

$$\sum_{i=1}^p V\{Y_i\} = \sum_{i=1}^p V\{z_i\} = p \quad \rho_{Y_i, Y_k} = \rho_{i,k} \sqrt{\lambda_i} \quad i, k=1, \dots, p$$

$(\lambda_i, \underline{e}_i) \equiv$ eigenvalue-eigenvector pairs of \underline{P}

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$$

Proportion of (standardized) population variance due

$$\text{to the } k^{\text{th}} \text{ component} \equiv \frac{\lambda_k}{p} \quad k=1, \dots, p$$

EXAMPLE - LPGA DATA

$$R = \begin{bmatrix} 1.000 & -.490 & .467 & -.179 \\ & 1.000 & .229 & .194 \\ & & 1.000 & .053 \\ & & & 1.000 \end{bmatrix}$$

$$\lambda_1 = 1.647 \quad \lambda_2 = 1.292 \quad \lambda_3 = 0.864 \quad \lambda_4 = 0.197$$

$$\underline{e}_1 = \begin{bmatrix} .715 \\ -.537 \\ .298 \\ -.334 \end{bmatrix} \quad \underline{e}_2 = \begin{bmatrix} -.205 \\ -.491 \\ -.775 \\ -.342 \end{bmatrix} \quad \underline{e}_3 = \begin{bmatrix} -.156 \\ .427 \\ .157 \\ -.877 \end{bmatrix} \quad \underline{e}_4 = \begin{bmatrix} -.650 \\ -.537 \\ .535 \\ -.050 \end{bmatrix}$$

# Components	1	2	3	4
Proportion of Variation	$\frac{1.647}{4} = .412$	$\frac{2.939}{4} = .735$	$\frac{3.803}{4} = .951$	1

$$r_{Y_1, z_1} = e_{11} \sqrt{\lambda_1} = .715 \sqrt{1.647} = .918 \quad r_{Y_1, z_2} = e_{12} \sqrt{\lambda_1} = -.537 \sqrt{1.647} = -.689$$

Principal Components for Special Covariance Matrices

$$1) \quad \Sigma = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \dots & \\ 0 & & & \sigma_{pp} \end{bmatrix} \quad \underline{e}_i = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\lambda \underline{e}_i = \sigma_{ii} \underline{e}_i \Rightarrow (\sigma_{ii}, \underline{e}_i) \equiv \text{eigenvalue/eigenvector pair}$$

$$R = I \quad \underline{e}_i = \underline{e}_i \quad (1, \underline{e}_i) \equiv \dots$$

2) Compound Symmetry $\Sigma = \begin{bmatrix} \sigma^2 & & & \\ & \sigma^2 & & \\ & & \dots & \\ \rho\sigma^2 & & & \sigma^2 \end{bmatrix}$ $\rho = \begin{bmatrix} 1 & & & \\ & \rho & & \\ & & \dots & \\ \rho & & & 1 \end{bmatrix}$

When $\rho > 0 \Rightarrow$

$$\lambda_1 = 1 + (p-1)\rho \quad e_1' = \left[\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}}, \dots, \frac{1}{\sqrt{p}} \right]$$

$$\lambda_2 = \lambda_3 = \dots = \lambda_p = 1 - \rho$$

$$e_2' = \left[\frac{1}{\sqrt{1 \times 2}}, \frac{-1}{\sqrt{1 \times 2}}, 0, \dots, 0 \right]$$

$$e_3' = \left[\frac{1}{\sqrt{2 \times 3}}, \frac{1}{\sqrt{2 \times 3}}, \frac{-2}{\sqrt{2 \times 3}}, 0, \dots, 0 \right]$$

$$e_i' = \left[\frac{1}{\sqrt{(i-1)i}}, \dots, \frac{1}{\sqrt{(i-1)i}}, \frac{-(i-1)}{\sqrt{(i-1)i}}, 0, \dots, 0 \right]$$

$$e_p' = \left[\frac{1}{\sqrt{(p-1)p}}, \dots, \frac{1}{\sqrt{(p-1)p}}, \frac{-(p-1)}{\sqrt{(p-1)p}} \right]$$

$$Y_1 = e_1' z = \frac{1}{\sqrt{p}} \sum_{i=1}^n z_i \quad \text{"equal weights" - index.}$$

$$\frac{\lambda_1}{p} = \frac{1 + (p-1)\rho}{p} = \rho + \frac{1-\rho}{p}$$

$z_i \sim N_p(0, \rho) \Rightarrow$ ellipsoids of constant density are cigar shaped w/ 1st principal

Component $Y_1 = \frac{1}{\sqrt{p}} \mathbf{1}' z$ (Projection of z onto

equiangular line $\mathbf{1}' = [1 \dots 1]$ Minor axes

occur in spherically symmetric directions perpendicular to major axis.

8.3 Summarizing Sample Variation by Principal Components

Observe $X_1, \dots, X_n \stackrel{iid}{\sim} (\mu, \Sigma) \Rightarrow \bar{X}, S, R$

$$a_i' \underline{x} = a_{i1} x_{j1} + \dots + a_{ip} x_{jp} \quad j=1, \dots, p$$

\Rightarrow sample mean $a_i' \bar{X}$, sample variance $a_i' S a_i$

$$\text{Cov}\{a_i' X_j, a_k' X_j\} = a_i' S a_k$$

Sample Principal Components and variances obtained in same manner as population, based on S , not Σ

$S = \{S_{ik}\} \equiv p \times p$ Sample variance-covariance matrix w/
eigenvalue-eigenvector pairs $(\hat{\lambda}_1, \hat{e}_1), \dots, (\hat{\lambda}_p, \hat{e}_p)$

$$\hat{Y}_i = \hat{e}_i' \underline{x} = \hat{e}_{i1} x_1 + \hat{e}_{i2} x_2 + \dots + \hat{e}_{ip} x_p \quad i=1, \dots, p$$

$$\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0 \quad x \text{ is any observation on } X_1, \dots, X_p$$

$$\text{Sample var } \{\hat{Y}_k\} = \hat{\lambda}_k \quad k=1, \dots, p$$

$$\text{Sample cov } \{\hat{Y}_i, \hat{Y}_k\} = 0 \quad i \neq k$$

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

$$r_{\hat{Y}_i, X_k} = \frac{\hat{e}_{ik} \sqrt{\hat{\lambda}_i}}{\sqrt{S_{kk}}} \quad i, k = 1, \dots, p$$

Sample Principal Components can be obtained from $\hat{\Sigma}, S, R$

Centered values: $\hat{y}_i = \hat{e}_i' (x - \bar{x})$ $i=1, \dots, p$ for any \underline{x}

Observed values: $\hat{y}_{ji} = \hat{e}_i' (x_j - \bar{x})$ $j=1, \dots, n$

$$\Rightarrow \bar{y}_i = \frac{1}{n} \sum_{j=1}^n \hat{e}_i' (x_j - \bar{x}) = \frac{1}{n} \hat{e}_i' \sum_{j=1}^n (x_j - \bar{x}) = 0$$

Example - Metals (Cd, Hg, Pb) in Mediterranean allscorers.

$n=34$

$$R = \begin{bmatrix} 1.0000 & 0.3112 & -0.0395 \\ & 1.0000 & -0.1065 \\ & & 1.0000 \end{bmatrix}$$

$$\lambda_1 = 1.3426$$

$$\lambda_2 = 0.9760$$

$$\lambda_3 = 0.6814$$

$$\underline{e}_1 = \begin{bmatrix} .6614 \\ .6912 \\ -.2911 \end{bmatrix}$$

$$\underline{e}_2 = \begin{bmatrix} .3157 \\ .0955 \\ .9441 \end{bmatrix}$$

$$\underline{e}_3 = \begin{bmatrix} .6804 \\ -.7163 \\ -.1550 \end{bmatrix}$$

Cumulative
Proportion: $\frac{1.3426}{3} = .4475$

$\frac{2.3186}{3} = .7729$

1

z_k	$\sqrt{\lambda_{1,z_k}}$	
z_1	.6614 $\sqrt{1.3426}$	= .7664
z_2	.6912 $\sqrt{1.3426}$	= .8609
z_3	-.2911 $\sqrt{1.3426}$	= -.3373

$\sqrt{\lambda_{2,z_k}}$	
.3157 $\sqrt{.9760}$	= .3119
.0955 $\sqrt{.9760}$	= .0943
.9441 $\sqrt{.9760}$	= .9327

$\sqrt{\lambda_{3,z_k}}$	
.6804 $\sqrt{.6814}$	= .5616
-.7163 $\sqrt{.6814}$	= -.5913
-.1550 $\sqrt{.6814}$	= -.1279

Interpretation

Mean of
Cd, Hg

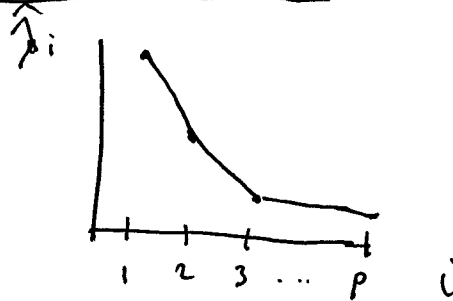
Pb

Difference
between
Cd, Hg

Number of Principal Components

Scree plot

Look for elbow/bend in plot.



EXAMPLE 8.4 - PROGRAM posted on webpage

Interpretation of Sample Principal Components

$\hat{y}_i = \hat{e}_i' (\underline{x} - \bar{x})$ realizations of population $y_i = e_i' (\underline{x} - \underline{\mu})$

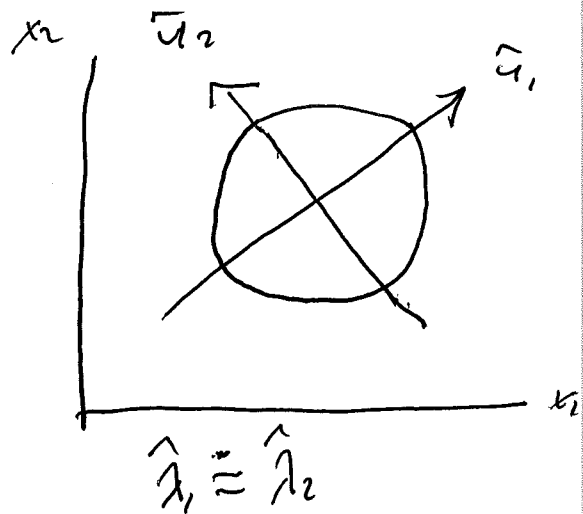
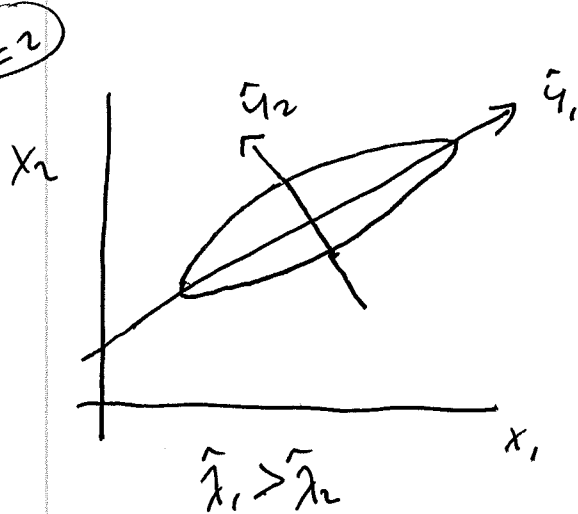
From sample values x_j , approximate $\underline{\mu}$ by \bar{x} , Σ by S

$S \equiv \text{pos. def.} \Rightarrow (\underline{x} - \bar{x})' S^{-1} (\underline{x} - \bar{x}) = c^2$ estimates

contours of $(\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu}) = c^2$

Len, μ s of hyperellipsoid axes $\approx \sqrt{\hat{\lambda}_i} \quad i=1, \dots, P$

$P=2$



Standardizing the Sample Principal Components

8.13

$$\tilde{z}_j = D^{-1/2} (x_j - \bar{x}) = \begin{bmatrix} \frac{x_{j1} - \bar{x}_1}{\sqrt{s_{11}}} \\ \vdots \\ \frac{x_{jp} - \bar{x}_p}{\sqrt{s_{pp}}} \end{bmatrix} \quad j=1, \dots, n$$

$$\underset{n \times p}{Z} = \begin{bmatrix} z_1' \\ z_2' \\ \vdots \\ z_n' \end{bmatrix} = \begin{bmatrix} \frac{x_{11} - \bar{x}_1}{\sqrt{s_{11}}} & \dots & \frac{x_{1p} - \bar{x}_p}{\sqrt{s_{pp}}} \\ \vdots & & \vdots \\ \frac{x_{n1} - \bar{x}_1}{\sqrt{s_{11}}} & \dots & \frac{x_{np} - \bar{x}_p}{\sqrt{s_{pp}}} \end{bmatrix}$$

$$\bar{z} = \frac{1}{n} Z' \underline{1}_n$$

$$\begin{aligned} S_z &= \frac{1}{n-1} \left(Z - \frac{1}{n} \underline{1} \underline{1}' Z \right)' \left(Z - \frac{1}{n} \underline{1} \underline{1}' Z \right) \\ &= \frac{1}{n-1} \left(Z - \underline{1} \bar{z}' \right)' \left(Z - \underline{1} \bar{z}' \right) = \frac{1}{n-1} Z' Z \end{aligned}$$

$$\begin{aligned} \text{Note: } \sum_{j=1}^n \left[\frac{(x_{ji} - \bar{x}_i)}{\sqrt{s_{ii}}} \frac{(x_{jk} - \bar{x}_k)}{\sqrt{s_{kk}}} \right] &= \frac{(n-1) s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}} \\ &= (n-1) r_{ik} \end{aligned}$$

$$\Rightarrow S_z = \frac{1}{n-1} Z' Z = R$$

z_1, \dots, z_n = standardized observations w/
covariance matrix R .

$$\hat{y}_i = \hat{e}_i' z = \hat{e}_{i1} z_1 + \dots + \hat{e}_{ip} z_p \quad i=1, \dots, p$$

where $(\hat{\lambda}_i, \hat{e}_i)$ is i^{th} eigenvalue-eigenvector pair of R

$$w/ \quad \hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$$

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

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

$$\text{Total (standardized) sample variance} = \text{tr}(R) = p = \hat{\lambda}_1 + \dots + \hat{\lambda}_p$$

$$r_{\hat{y}_i, z_k} = \hat{e}_{ik} \sqrt{\hat{\lambda}_i} \quad i, k=1, \dots, p$$

Proportion of standardized sample variance due to p^{th} PC:

$$\frac{\hat{\lambda}_i}{p} \quad i=1, \dots, p$$

EXAMPLE 8.5 - ON WEBPAGE

8.4 Graphing Principal Components

$$\underline{x}_j = \begin{pmatrix} x_{j1} \\ \vdots \\ x_{jp} \end{pmatrix} \hat{e}_1 + \dots + \begin{pmatrix} x_{j1} \\ \vdots \\ x_{jp} \end{pmatrix} \hat{e}_p = \hat{y}_{j1} \underline{e}_1 + \dots + \hat{y}_{jp} \underline{e}_p$$

1. Check for normality: (a) Scatterplots of first few PCs
(b) Q-Q plots of sample values for each PC.
2. Scatter diagrams and Q-Q plots for last few PCs to identify suspect data cases.

EXAMPLE 8.7 R PROGRAM ON WEBSITE

8.5 Large-Sample Inferences

Assumptions:

$\underline{x}_1, \dots, \underline{x}_n \equiv$ Random sample from Normal Population

Eigenvalues of Σ are distinctly positive $\lambda_1 > \dots > \lambda_p > 0$

(Case where the number of equal eigenvalues is known is an exception (e.g. Testing Compound Symmetry)).

Large-sample Theory of $\hat{\lambda}' = [\hat{\lambda}_1, \dots, \hat{\lambda}_p]$, $\hat{e}_1, \dots, \hat{e}_p$ for S

$$\textcircled{1} \Lambda = \text{diag} \{ \lambda_i \} \text{ for } \Sigma \Rightarrow \sqrt{n}(\hat{\lambda} - \lambda) \sim N_p(\underline{0}, 2\Lambda^2)$$

$$\textcircled{2} \text{Let } E_i = \lambda_i \sum_{\substack{k=1 \\ k \neq i}}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} \underline{e}_k \underline{e}_k' \Rightarrow \sqrt{n}(\hat{e}_i - \underline{e}_i) \sim N_p(\underline{0}, E_i)$$

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

$i=1, \dots, p$

$$\begin{aligned}
 \textcircled{1} &\Rightarrow P_r \left\{ -z\left(\frac{\alpha}{2}\right) \sqrt{2\lambda_i} \leq \sqrt{n}(\hat{\lambda}_i - \lambda_i) \leq z\left(\frac{\alpha}{2}\right) \sqrt{2\lambda_i} \right\} = 1 - \alpha \\
 &= P_r \left\{ -z\left(\frac{\alpha}{2}\right) \sqrt{2} \leq \sqrt{n} \left(\frac{\hat{\lambda}_i}{\lambda_i} - 1 \right) \leq z\left(\frac{\alpha}{2}\right) \sqrt{2} \right\} \\
 &= P_r \left\{ -z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}} \leq \frac{\hat{\lambda}_i}{\lambda_i} - 1 \leq z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}} \right\} \\
 &= P_r \left\{ 1 - z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}} \leq \frac{\hat{\lambda}_i}{\lambda_i} \leq 1 + z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}} \right\} \\
 &= P_r \left\{ \frac{1}{\lambda_i} \left(1 - z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}} \right) \leq \frac{1}{\lambda_i} \leq \frac{1}{\lambda_i} \left(1 + z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}} \right) \right\} \\
 &= P_r \left\{ \frac{\hat{\lambda}_i}{1 + z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}}} \leq \lambda_i \leq \frac{\hat{\lambda}_i}{1 - z\left(\frac{\alpha}{2}\right) \sqrt{\frac{2}{n}}} \right\}
 \end{aligned}$$

\equiv Large-sample $100(1-\alpha)\%$ CI for λ_i
 Bonferroni simultaneous CIs can be computed
 for m λ_i 's by $z\left(\frac{\alpha}{2}\right) \rightarrow z\left(\frac{\alpha}{2m}\right)$

$\textcircled{2}$ for E_i , must use estimates $\{\hat{\lambda}_i\}, \{\hat{e}_i\}$ to
 obtain estimated standard errors.

Testing for equal Correlation Structure

$$H_0: \rho = \rho_0 = \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \dots & 1 \end{bmatrix} \quad H_1: \rho \neq \rho_0$$

Lewley's procedure

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

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

Reject H_0 in favor of H_A if:

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