

### III Independent Samples

See Also Section 4.3 of Khattree & Naik

Suppose  $n$  units are randomly assigned to either treatment 1 or treatment 2.

Randomization is essential to make the inference mathematically valid. Although it does not guarantee elimination of bias, one hopes it will reduce any systematic bias due to inhomogeneity among the units with respect to uncontrolled variables.

- Assume
1.  $X_{11}, \dots, X_{1n_1}$  IID  $(\mu_1, \Sigma_1)$
  2.  $X_{21}, \dots, X_{2n_2}$  IID  $(\mu_2, \Sigma_2)$
  3. the two samples are independent

For small  $n_1, n_2$  further assume

- A a) the two populations are normal
- b)  $\Sigma_1 = \Sigma_2 (= \Sigma)$

# A Small sample inference

## Estimation

$$\hat{\mu}_1 = \bar{x}_1, \quad \hat{\mu}_2 = \bar{x}_2, \quad \hat{\Sigma} = S_p \stackrel{\text{def}}{=} \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1+n_2-2}$$

$$\mathbb{E} [\bar{x}_1 - \bar{x}_2] = \hat{\mu}_1 - \hat{\mu}_2 = \frac{1}{n_1+n_2-2} \sum_{k=1}^2 \sum_{j=1}^{n_k} (x_{kj} - \bar{x}_k)(x_{kj} - \bar{x}_k)'$$

$$\text{Cov} [\bar{x}_1 - \bar{x}_2] = \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \hat{\Sigma}$$

$$\text{Cov} [\bar{x}_1 - \bar{x}_2] = \left( \frac{1}{n_1} + \frac{1}{n_2} \right) S_p$$

## Testing

$$H_0: \mu_1 - \mu_2 = \delta_0 \quad \text{vs} \quad H_1: \mu_1 - \mu_2 \neq \delta_0$$

$$T^2 = (\bar{x}_1 - \bar{x}_2 - \delta_0)' \left[ \left( \frac{1}{n_1} + \frac{1}{n_2} \right) S_p \right]^{-1} (\bar{x}_1 - \bar{x}_2 - \delta_0)$$

$$\sim \frac{(n_1+n_2-2)p}{n_1+n_2-1-p} F_{p, n_1+n_2-p-1}$$

since when  $H_0$  holds,

indep

$$\begin{cases} \bar{x}_1 - \bar{x}_2 - \delta_0 \sim N_p(\mathbf{0}, (\frac{1}{n_1} + \frac{1}{n_2}) \hat{\Sigma}) \\ S_p(n_1+n_2-2) \sim W_{n_1+n_2-2}(\hat{\Sigma}) \end{cases}$$

See Result on page 33

(right after the comparison of univariate and multivariate)

If  $U \sim N_p(\mathbf{0}, \hat{\Sigma})$ , independent

$A \sim W_m(\hat{\Sigma})$

$$\text{then } U' A^{-1} U \sim \frac{p}{m-p+1} F_{p, m-p+1}$$

Confidence Region  $(1-\alpha)100\%$  C.R. for  $\delta = \mu_1 - \mu_2$

$$\{ \delta : [\delta - (\bar{x}_1 - \bar{x}_2)]' \left[ \left( \frac{1}{n_1} + \frac{1}{n_2} \right) S_p \right]^{-1} [\delta - \bar{x}_1 - \bar{x}_2]$$

$$\leq \frac{(n_1+n_2-2)p}{n_1+n_2-1-p} F_{p, n_1+n_2-1-p}$$

= ellipsoid with center  $\bar{x}_1 - \bar{x}_2$  and axes

$$\sqrt{\frac{(n_1+n_2-2)p}{n_1+n_2-1-p} \left( \frac{1}{n_1} + \frac{1}{n_2} \right) F_c} \lambda_i e_i$$

where  $F_c = F_{p, n_1+n_2-1-p}(\alpha)$

where  $S_p = \sum_{i=1}^p \lambda_i e_i e_i'$  is the standard spectral decomposition

Simultaneous C.I. for  $\ell' \delta$

$$\ell' (\bar{x}_1 - \bar{x}_2) \pm \sqrt{\frac{(n_1+n_2-2)p}{n_1+n_2-1-p} F_{p, n_1+n_2-1-p}(\alpha)} \sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) \ell' S_p \ell}$$

In particular, C.I. for  $\mu_{1i} - \mu_{2i}$

$$(\bar{x}_{1i} - \bar{x}_{2i}) \pm \sqrt{\frac{(n_1+n_2-2)p}{n_1+n_2-1-p} F_{p, n_1+n_2-1-p}(\alpha)} \sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) (S_p)_{ii}}$$

Bonferroni C.I. for  $\mu_{1i} - \mu_{2i}$

$$(\bar{x}_{1i} - \bar{x}_{2i}) \pm t_{\frac{n_1+n_2-2}{n_1+n_2-1-p}} \left( \frac{\alpha}{2p} \right) \sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) (S_p)_{ii}}$$

## B Large Sample inference ( $\Sigma_1 = \Sigma_2$ )

1 Robust wrt non-normality

2 Replace  $\frac{(n_1+n_2-2)p}{n_1+n_2-1-p} F_{p, n_1+n_2-1-p}$

by  $\chi_p^2$

and follow recipe A. That is,

$$T^2 = (\bar{x}_1 - \bar{x}_2 - \delta_0)' \left[ \left( \frac{1}{n_1} + \frac{1}{n_2} \right) S_p \right]^{-1} (\bar{x}_1 - \bar{x}_2 - \delta_0) \stackrel{\text{approx}}{\sim} \chi_p^2$$

$$\text{CR } \{ \delta : [\delta - (\bar{x}_1 - \bar{x}_2)]' \left[ \left( \frac{1}{n_1} + \frac{1}{n_2} \right) S_p \right]^{-1} [\delta - (\bar{x}_1 - \bar{x}_2)] \leq \chi_p^2(\alpha) \}$$

$$\text{SCI for } \delta : \delta' (\bar{x}_1 - \bar{x}_2) \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) \delta' S_p \delta}$$

$$\text{for } \mu_{1i} - \mu_{2i} : (\bar{x}_{1i} - \bar{x}_{2i}) \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) (S_p)_{ii}}$$

$$\text{Bon.C.I for } \mu_{1i} - \mu_{2i} \text{ is } (\bar{x}_{1i} - \bar{x}_{2i}) \pm t_{n_1+n_2-2} \left( \frac{\chi_p^2}{n_1+n_2-2} \right) \sqrt{\left( \frac{1}{n_1} + \frac{1}{n_2} \right) S_{p,ii}}$$

C If  $\Sigma_1 \neq \Sigma_2$  and  $n_1, n_2$  are moderate,

Try some transformation on variables that have drastically different variance estimates than others.

# D Large Sample Inference $\Sigma_1 \neq \Sigma_2$

Estimation

$$\hat{\mu}_1 = \bar{x}_1, \quad \hat{\mu}_2 = \bar{x}_2, \quad \hat{\Sigma}_1 = S_1, \quad \hat{\Sigma}_2 = S_2$$

$$(\bar{x}_1 - \bar{x}_2 - \delta_0)' \left( \frac{S_1}{n_1} + \frac{S_2}{n_2} \right)^{-1} (\bar{x}_1 - \bar{x}_2 - \delta_0) \stackrel{\text{approx}}{\sim} \chi_p^2$$

Since  $\bar{x}_1 - \bar{x}_2 \sim N(\delta, \frac{\Sigma_1}{n_1} + \frac{\Sigma_2}{n_2})$ , under H<sub>0</sub>.

(1-a)100% CR for  $\delta = \hat{\mu}_1 - \hat{\mu}_2$  is

$$\left\{ \delta: [\delta - (\bar{x}_1 - \bar{x}_2)]' \left( \frac{S_1}{n_1} + \frac{S_2}{n_2} \right)^{-1} [\delta - (\bar{x}_1 - \bar{x}_2)] \leq \chi_p^2(\alpha) \right\}$$

= ellipsoid with center  $\bar{x}_1 - \bar{x}_2$  and axes

$$\sqrt{\chi_p^2(\alpha)} e_i \quad \text{where} \quad \frac{S_1}{n_1} + \frac{S_2}{n_2} = \sum_{i=1}^p \gamma_i e_i e_i'$$

(1-a)100% SCI for  $\ell' \delta$  is

$$\ell' (\bar{x}_1 - \bar{x}_2) \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\frac{\ell' S_1 \ell}{n_1} + \frac{\ell' S_2 \ell}{n_2}}$$

In particular, SCI for  $\hat{\mu}_{1i} - \hat{\mu}_{2i}$  is

$$(\bar{x}_{1i} - \bar{x}_{2i}) \pm \sqrt{\chi_p^2(\alpha)} \sqrt{\frac{S_{1ii}}{n_1} + \frac{S_{2ii}}{n_2}}$$

# Comparison of multivariate population means.

Suppose ①  $\underline{X}_{l1}, \dots, \underline{X}_{ln_l}$  is a r.s from  $N_p(\mu_l, \Sigma)$   
 $l=1, \dots, g$

② The  $g$  samples are independent

Note: The population covariance matrices are same.

$H_0: \mu_1 = \dots = \mu_g$  or equivalently

$H_0: \tau_1 = \dots = \tau_g$  where  $\tau_l = \mu_l - \mu$  additional effect

One Way

MANOVA model

$$\underline{X}_{lj} = \underbrace{\mu}_{\text{overall mean}} + \underbrace{\tau_l}_{\substack{\text{treatment effect} \\ l^{\text{th}} \text{ treatment}}} + \underbrace{\varepsilon_{lj}}_{\substack{\text{error term} \\ \text{IID } N_p(0, \Sigma)}} \quad (\mu + \mu_l - \mu = \frac{1}{g} \sum_l \mu_l)$$

$$\text{Similarly, } \sum_{l=1}^g \tau_l = 0$$

Source	df	SS
treatment (Between) (Treatment) and cross products	$g-1$	$B = SST_r = \sum_{l=1}^g n_l (\bar{X}_l - \bar{X}) (\bar{X}_l - \bar{X})'$
residual (Within) sum of squares (Residual/Error) and cross products	$n-g$	$W = SS_{Res} = \sum_{l=1}^g \sum_{j=1}^{n_l} (X_{lj} - \bar{X}_l) (X_{lj} - \bar{X}_l)'$
Total (Corrected)	$n-1$	$T = SST_{Total} = \sum_{l=1}^g \sum_{j=1}^{n_l} (X_{lj} - \bar{X}) (\bar{X}_l - \bar{X})'$

$$n = \sum_{l=1}^g n_l, \quad T = B + W \quad X_{lj} - \bar{X} = \underbrace{X_{lj} - \bar{X}_l}_{\text{Within treatment}} + \underbrace{\bar{X}_l - \bar{X}}_{\text{between treatment}}$$

Wilks Lambda (ratio of generalized variance under full model and under  $H_0$ )

$$\Lambda^* = \frac{|\hat{\Sigma}|}{|\hat{\Sigma}_0|} = \frac{|W|}{|B + W|} = \frac{S}{\sum_{i=1}^n \left( \frac{1}{\lambda_i} \right)} \quad \text{where } (\lambda_i) \text{ are eigen values of } W^{-1}$$

Note  $W \sim W_{n-q}(\Sigma)$

$B \sim W_{q-1}(\Sigma)$

Distribution of  $\Lambda^*$  under  $H_0$  provided  $H_0$  is true

$$\text{large } n, \quad -\left(n-1-\frac{p+q}{2}\right) \ln \Lambda^* \text{ approx } \chi^2_{p(q-1)}$$

Since  $p > q = \# \text{ parameters specified by } H_0 = p(q-1)$

Moderate  $n$  (special cases)

$$\frac{1 - \Lambda^{*\frac{1}{p}}}{\Lambda^{*\frac{1}{p}}} \sim \frac{p(q-1)}{p(n-q-2)(p-1)} F_{p(q-1)}, p(n-q-2(p-1)}$$

$$q = 2, 3 \quad \frac{1 - \Lambda^{*\frac{1}{q-1}}}{\Lambda^{*\frac{1}{q-1}}} \sim \frac{p}{(n-q)-(p-1)} F_{p(q-1)}, [n-q-p+1](q-1)$$

Read 6.7 for two-way Multivariate analyses of Variance.

Example

			Sample Average
8	10	9	$\left(\begin{matrix} 9 \\ 8 \end{matrix}\right) = \bar{x}_1$
5	9	10	
3	5		$\left(\begin{matrix} 4 \\ 4 \end{matrix}\right) = \bar{x}_2$
3	5		
4	9	8	$\left(\begin{matrix} 7 \\ 8 \end{matrix}\right) = \bar{x}_3$
7	9	9	
			$\left(\begin{matrix} 7 \\ 7 \end{matrix}\right) = \bar{x}$

$$g=3, n=8$$

$$\Lambda^* = -\frac{|W|}{|B+W|} = -\frac{224}{860} = -2.604651$$

$$\begin{pmatrix} 3 & 1 & -3 \\ 3 & 9 & 9 \\ 2 & 9 & 9 \end{pmatrix} \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} 3 \begin{pmatrix} 4 & 2 \\ 2 & 1 \end{pmatrix} \\ 3 \begin{pmatrix} 0 & 10 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$$

$$B = \begin{bmatrix} +18 & +6 \\ +6 & +3 \\ +18 & +3 \end{bmatrix} = \begin{bmatrix} 30 & 24 \\ 24 & 24 \end{bmatrix}$$

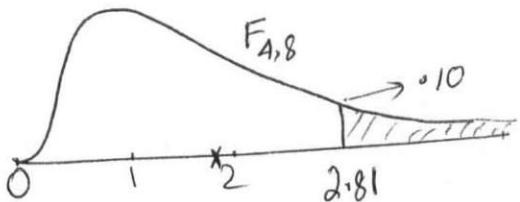
$$W = \begin{bmatrix} 111 & 31 \\ 111 & 11 \\ 9111 & 31 \\ 31 & 914 \\ 11 & 11 \\ 3 & 1 \end{bmatrix} = \begin{bmatrix} 18 & 10 \\ 10 & 18 \end{bmatrix}$$

$$B+W = \begin{bmatrix} 48 & 34 \\ 34 & 42 \end{bmatrix} \\ \begin{pmatrix} 1 & 10 \\ -3 & 12 \end{pmatrix} \begin{pmatrix} -1 & -3 \\ 1 & 1 \\ 0 & 2 \end{pmatrix}$$

$$\frac{H(\Lambda^*)^{\frac{1}{p}}}{(\Lambda^*)^{\frac{1}{q}}} \sim \frac{P}{n-q-p+1} F_{P(q-1), (n-q-p+1)(q-1)} \\ = \frac{2}{4} F_{4,8}$$

$$F_{\text{obs}} = \frac{4}{2} \frac{1 - \sqrt{\Lambda^*}}{\sqrt{\Lambda^*}} = 2 (0.9594095) = 1.9188$$

$$F_{4,8}(0.10) = 2.81 \quad \text{p-value} > 0.10$$



Simultaneous C.I. for treatment effects

$\binom{g}{2}$  possible comparison (pairwise)

use  $\alpha' = \frac{\alpha}{P(\binom{g}{2})}$  for each test

$H_0: \bar{x}_{ki} = \bar{x}_{li}; i \leq k \neq l \leq g, i \leq l$  fixed

$$\widehat{x}_{ki} - \widehat{x}_{li} = \bar{x}_{ki} - \bar{x}_{li} \sim N(0, (\frac{1}{n_k} + \frac{1}{n_l}) \sigma_{ii}^2)$$

$$\widehat{\sigma}_{ii} = \frac{\omega_{ii}}{n-g} \quad \text{where} \quad W = \sum_{k=1}^g \sum_{j=1}^{n_k} (x_{kj} - \bar{x}_k)(x_{kj} - \bar{x}_k)'$$

(1- $\alpha$ ) 100% S.C.I. for  $\bar{x}_{ki} - \bar{x}_{li}$  ( $k \neq l$ ) are

$$(\bar{x}_{ki} - \bar{x}_{li}) \pm t_{n-g} \left( \frac{\alpha}{P(g(g-1))} \right) \sqrt{\left( \frac{1}{n_k} + \frac{1}{n_l} \right) \frac{\omega_{ii}}{n-g}}$$

```

options ls=80 ps=47 nodate nonumber; code6-2
title1 h=2 'Hook-billed kites';
title2 h=1 'JW: T5-12 T6-11';
data length1;
infile 'Z:\My Documents\Teaching\Stat524\Fall 2010\Data set\T5-12.dat' firstobs=1;
input tail wing;
sex='female';
run;

data length2;
infile 'Z:\My Documents\Teaching\Stat524\Fall 2010\Data set\T6-11.dat' firstobs=1;
input tail wing;
sex='male';
if tail >= 284 then delete;
run;

data length;
set length1 length2;
run;

proc print data=length;
run;

proc corr data=length cov;
var tail wing;
by sex;
run;
/*----- sex=female -----*/
2 Variables: tail wing
Covariance Matrix, DF = 44
tail          tail          wing
tail          120.694949      122.3459596
wing          122.3459596     208.5404040
----- sex=male -----
2 Variables: tail wing
Covariance Matrix, DF = 43
tail          tail          wing
tail          88.3694503     90.0073996
wing          90.0073996     172.8350951
----- simple statistics -----
Simple Statistics
Variable N Mean Std Dev Sum Minimum Maximum
tail    45 193.62222 10.98613 8713 173.00000 216.00000
wing   45 279.77778 14.44093 12590 245.00000 310.00000
----- simple statistics -----
Simple Statistics
Variable N Mean Std Dev Sum Minimum Maximum
tail    44 187.15909 9.40050 8235 170.00000 212.00000
wing   44 280.95455 13.14668 12362 254.00000 310.00000
*/
proc iml;
nf=45;
p=2;
avef={193.62222, 279.77778};
Sf={120.694945 122.3459596;
122.3459596 208.5404040};

```

## code6-2

```

nm=44;
avem={187.15909, 280.95455};
Sm={88.3694503 90.0073996,
     90.0073996 172.8350951};
diff=avef-avem;
n=nf+nm;
df1=p;
df2=n-p-1;
Sp=((nf-1)*Sf+(nm-1)*Sm)/(n-2);
Spi=inv((1/nf+1/nm)*Sp);
T2=diff*Spi*diff;
alpha=0.05;
falpha=finv(1-alpha,df1,df2);
f=(n-2)*df1/df2*falpha;
print diff Sp, T2 f alpha df1 df2;
c1=sqrt(f);
c2=tinv(1-alpha/(2*p),n-2);
simulcl=diff-c1*sqrt((1/nf+1/nm)*vecdiag(Sp));
simuucl=diff+c1*sqrt((1/nf+1/nm)*vecdiag(Sp));
bonflcl=diff-c2*sqrt((1/nf+1/nm)*vecdiag(Sp));
bonfuc1=diff+c2*sqrt((1/nf+1/nm)*vecdiag(Sp));
print c1 c2, simulcl simuucl, bonflcl bonfuc1;
run;
/*
          diff      Sp
          6.46313 104.71798 106.36253
          -1.17677 106.36253 190.89295

```

T2	f	alpha	df1	df2
24.964901	6.2772565	0.05	2	86
c1	c2			
2.5054454	2.2808563			
simulcl	simuucl			
1.0273988	11.898861			
-8.515861	6.1623214			
bonflcl	bonfuc1			
1.5146599	11.4116			
-7.857982	5.5044423			

```

*/
proc GLM data=length;
class sex;
model tail wing =sex;
manova h=sex/printe printh;
means sex/bon alpha=0.025 lines;
means sex/bon alpha=0.025 cldiff;
run;
quit;
-----
```

## Assignment #5

The day assigned: Thursday, October 14, 2010  
 The day due: Tuesday, October 26, 2010

1. Read Text: Chapter 6
2. Assigned problems: Exercises: 6.5, 6.8

# Lecture 16-17

## Chapter 8 Analysis of covariance structure principal Component Analysis

### 8.1 Introduction

#### (a) what is PCA?

- <1> PCA is a method for re-expressing multivariate data. It allows to reorient the data so that the first few dimensions account for as much of the available information as possible.
- <2> If there is substantial redundancy present in the data set, then it may be possible to account for most of the information in the original data set with a relatively small number of dimensions.

#### (b) Objectives

- <1> Dimension (data) reduction: When there are many dimensions (variables), it is difficult to comprehend or even visualize the pattern of association among them.

When we examine relationship among  $p$  (very large) variables, we often look at the covariance matrix.

Sometimes, it is possible to explain the structure with only  $K$  ( $K < p$ ) linear combinations of the original  $p$  variables without loss of much information.

$n$  observation on  $p$  variables reduced to  
 $n$  observation on  $k$  ( $k \leq p$ ) variables

- (2) identifying patterns of association among variables, revealing relationship and the allowing interpretations that would not ordinarily result.
- (3) serving as intermediate steps in regression, clustering factor analysis, etc.

## 8.2 population principal components

### (a) principal components

(1) our objective is to find linear combinations of the original variable  $X = (X_1, \dots, X_p)$  which are uncorrelated with variance as large as possible

(2) Find linear combinations of  $p$  random variables

$$X_1, \dots, X_p$$

$$Y_1 = a_1' X = a_{11} X_1 + \dots + a_{1p} X_p$$

$$Y_2 = a_2' X = a_{21} X_1 + \dots + a_{2p} X_p$$

⋮

$$Y_p = a_p' X = a_{p1} X_1 + \dots + a_{pp} X_p$$

such that

$$Y_1, \dots, Y_p$$

- are uncorrelated
  - represent a new coordinate system obtained by rotating the original one with  $x_1, \dots, x_p$  as axes
  - represent the directions with maximum variability

(3) Let  $\text{cov}(X) = \Sigma$ , then

- $$\bullet \text{Var}(Y_c) = \text{Var}(a_c' X) = a_c' \Sigma a_c \quad c=1, \dots, p$$

$$\text{COV}(Y_i, Y_k) = \alpha_i' \Sigma \alpha_k = 0, \quad i \neq k = 1, \dots, p$$

- The 1<sup>st</sup> principal component =  $a_1' x$  that maximizes  $\text{var}(a' x)$  subject to  $a'a = 1$

$$\max_a \frac{a' \Sigma q}{a'a} = \lambda_1$$

$a = e_1$

The 2<sup>nd</sup>    - - - - -     $a_2'x$     ( $a_2' \neq a_2 = 1$ )

$$\max_{a \perp e_1} \frac{a' \Sigma a}{a'a} = \lambda_2$$

$a = e_2$

And  $\text{cor}(c_1 x, c_2 x) = 0$  (a.e.)

$$\text{The } i\text{th} \quad \dots = a_i x \quad (a_i' a_i = 1)$$

... - - - - - and  $\text{cov}(a_i'x, a_j'x) = 0$ ,  $\forall 1 \leq j < i$

(4) Let  $\Sigma = \text{cov}(X)$  have eigenvalue-eigen vector pairs

$(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$  with

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$  and  $\|e_i\| = 1$  for  $i \in p$

Then

- $$Y_i = e_i' x = e_{i1} x_1 + \dots + e_{ip} x_p \quad i=1, \dots, p$$

- $$\bullet \text{Var}(Y_i) = e_i' \Sigma e_i = \lambda_i$$

$$\text{Cor}(e_i, e_j) = 0 \quad i \neq j \quad (e_i' \Sigma e_j = \cancel{e_i' \lambda_j e_j} = 0)$$

$$\sum_{c=1}^p \text{Var}(Y_c) = \sum_{c=1}^p \lambda_c = \text{trace}(\Sigma) = \sum_{c=1}^p \text{Var}(X_c) = \sum_{c=1}^p \sigma_{cc}$$

↑

Total variance due  
to principal components      Total population  
variance

(5) Total variation explained by principal component

- $\frac{\lambda_i}{\sum \lambda_i} = \text{proportion of total variance explained by the } i\text{th principal component}$

thus, we can replace variables by the first a few principal components if they can account for most variation

(6) Importance of  $X_k$  to  $Y_i$

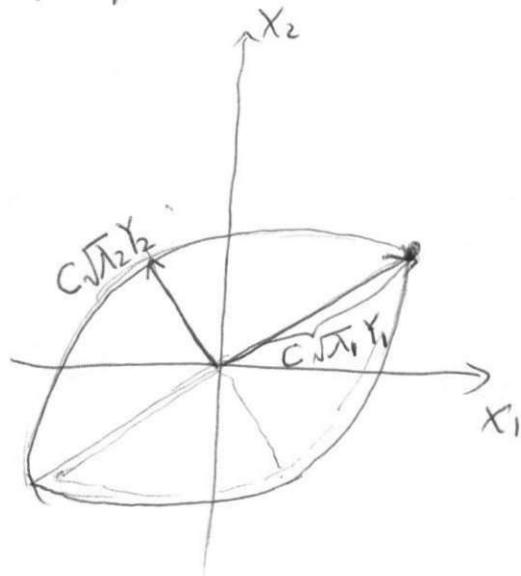
- $e_{ik} \text{ in } e_i' = (e_{i1}, \dots, e_{ik}, \dots, e_{ip})$  measures the importance of the  $k$ th variable to the  $i$ th principal component

$$Y_i = e_{i1}X_1 + \dots + e_{ik}X_k + \dots + e_{ip}X_p$$

$$\rho_{Y_i, X_k} = \frac{\text{cov}(Y_i, X_k)}{\sqrt{\text{Var}(Y_i)} \sqrt{\text{Var}(X_k)}} = \frac{\lambda_i e_{ik}}{\sqrt{\lambda_i} \sqrt{\sigma_{kk}}} = \frac{\sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} e_{ik}$$

(7) If  $X \sim N_p(\mu, \Sigma)$ , then contours of constant density  $(X-\mu)' \Sigma^{-1} (X-\mu) = c^2$  are ellipsoids centered at  $\mu$  with  $\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_p}$  as the ~~half~~ axes.

For  $p=2$  and  $\mu=0$



Example (8.1)

Suppose  $V(X) = \Sigma$ ,  $X = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$ ,  $\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$

- eigenvalue - eigen vector pairs

$$\lambda_1 = 5.83 \quad e_1' = (0.383, -0.924, 0)$$

$$\lambda_2 = 2.00 \quad e_2' = (0, 0, 1)$$

$$\lambda_3 = 0.17 \quad e_3' = (0.924, 0.383, 0)$$

$$(\lambda_1 + \lambda_2 + \lambda_3) = \sum \sigma_{ii} = 1 + 5 + 2 = 8$$

- the principle components

$$Y_1 = e_1' X = 0.383 X_1 - 0.924 X_2$$

$$Y_2 = e_2' X = X_3$$

$$Y_3 = e_3' X = 0.924 X_1 + 0.383 X_2$$

- proportion of total variance accounted by principal components

$$Y_1 : \frac{\lambda_1}{\sum \lambda_i} = \frac{5.83}{8} = 73\% \quad \left. \begin{array}{l} \\ \end{array} \right\} 98\%$$

$$Y_2 : \frac{\lambda_2}{\sum \lambda_i} = \frac{2}{8} = 25\%$$

$x_1, x_2, x_3$  can be replaced by  $Y_1$  and  $Y_2$

- The importance of  $x_i$  to  $Y_1$

$$Y_1 : |e_{12}| = |-0.924| > |e_{ii}| = |0.383|$$

$$P_{Y_1, X_2} = \sqrt{\frac{\lambda_1}{\sigma_{22}}} e_{12} = -0.998$$

$$P_{Y_1, X_1} = \sqrt{\frac{\lambda_1}{\sigma_{11}}} e_{ii} = \sqrt{\frac{5.83}{1}} 0.383 = 0.925$$

$x_1$  and  $x_2$  are almost equally important to  $Y_1$

with  $x_2$  being slightly more important.

- (b) principal components from standardized variable  
(correlation matrix)

<1> standardized variable and correlation matrix

- Let  $EX = (\mu_1, \dots, \mu_p)', \text{COV}(X) = \Sigma = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1p} \\ \vdots & \ddots & \vdots \\ \sigma_{p1} & \cdots & \sigma_{pp} \end{pmatrix}$

- Let  $Z_1 = \frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} , \dots , Z_p = \frac{x_p - \mu_p}{\sqrt{\sigma_{pp}}}$

$$\text{Then } Z = \begin{pmatrix} Z_1 \\ \vdots \\ Z_p \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{\sigma_{11}}} & 0 & \cdots & 0 \\ 0 & \frac{1}{\sqrt{\sigma_{22}}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\sqrt{\sigma_{pp}}} \end{pmatrix} \begin{pmatrix} X_1 - \mu_1 \\ \vdots \\ X_p - \mu_p \end{pmatrix} = (\text{diag}(\Sigma)^{-\frac{1}{2}})(X - \mu)$$

$$E(Z) = 0$$

$$\text{cov}(Z) = \begin{pmatrix} 1 & p_{12} & \cdots & p_{1p} \\ p_{21} & 1 & \cdots & p_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ p_{p1} & p_{p2} & \cdots & 1 \end{pmatrix} = P$$

(2) Let  $(\lambda_1, e_1), \dots, (\lambda_p, e_p)$  be the eigenvalue-eigenvector pair of  $P$  then

- $Y_i = e_i' Z = e_i' (\text{diag} \Sigma)^{-\frac{1}{2}} (X - \mu)$
- $\sum \text{var}(Y_i) = \sum \lambda_i = p = \sum \text{var}(Z_i)$
- $P_{Y_i, Z_k} = \frac{\sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} e_{ik} = \sqrt{\lambda_i} e_{ik} \quad i=1, \dots, p \quad k=1, \dots, p$

proportion of total variance explained by

$$Y_i : \frac{\lambda_i}{\sum \lambda_i} = \frac{\lambda_i}{p}$$

(3) Example (8.2)

Consider  $X_1$  and  $X_2$  with  $\Sigma = \begin{bmatrix} 1 & 4 \\ 4 & 100 \end{bmatrix}$

$$\text{The } P = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 1 \end{bmatrix}$$

For  $\Sigma$  we have

$$\lambda_1 = 100.16 \quad e_1 = [0.040, 0.999] \quad ; \quad \frac{\lambda_1}{\lambda_1 + \lambda_2} = 99.2\%$$

$$\lambda_2 = 0.84 \quad e_2 = [0.999, -0.040] \quad ; \quad \frac{\lambda_2}{\lambda_1 + \lambda_2} = 0.8\%$$

$$Y_i = e_i' X \quad i=1,2$$

$$e_{12} = 0.999 > e_{11} = 0.04$$

$$P_{Y_1, X_2} = \frac{e_{12}\sqrt{\lambda_1}}{\sqrt{\sigma_{22}}} = \frac{0.999\sqrt{100.16}}{\sqrt{100}} = 0.9998$$

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

For P we have

$$\lambda_1 = 1.4 \quad e_1' = (0.707, 0.707)$$

$$\lambda_2 = 0.6 \quad e_2' = (0.707, -0.707)$$

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = 70\%, \quad \frac{\lambda_2}{\lambda_1 + \lambda_2} = 30\%$$

$$e_{11} = e_{12}$$

$$P_{Y_1, X_1} = \frac{e_{11}\sqrt{\lambda_1}}{\sqrt{P_{11}}} = \frac{e_{12}\sqrt{\lambda_1}}{\sqrt{P_{22}}} = P_{Y_1, X_2} = 0.837$$

In this case, we should first do standardization before applying PCA.

(4) When to standardization

- standardization is not consequential and can make difference

- standardize if

- 1°) Variables are in widely different units

- 2°) Variable are measured in scales with widely different ranges

- If the response are reasonably commensurable, then covariance form has a greater statistical appeal.

c) principal components for special structured covariance matrix

$$(1) \text{ Diagonal matrix } \Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ & \ddots & & \\ & & \sigma_p^2 & \\ 0 & & & \end{bmatrix}$$

- $\lambda_i = \sigma_i^2 \quad e_i' = (0, \dots, 0, \underset{i\text{th}}{1}, 0, \dots, 0)$

$$Y_i = e_i' X = X_i$$

(No need for Principal components)

- Intra-class covariance (Equicorrelation) structure

$$\Sigma = \begin{bmatrix} \sigma^2 & \rho\sigma^2 & \cdots & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 & \cdots & \rho\sigma^2 \\ \vdots & \vdots & \ddots & \vdots \\ \rho\sigma^2 & \rho\sigma^2 & \cdots & \sigma^2 \end{bmatrix}$$

$$\rho = \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix}$$

$$= \frac{\Sigma}{\sigma^2} (1-\rho) I + \rho \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$$

For  $\rho$ .

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

$$\lambda_2 = 1-\rho \quad e_2' = \frac{1}{\sqrt{1-\rho}} (1, -1, \dots, 0)$$

$$\lambda_p = 1-\rho \quad e_p' = \frac{1}{\sqrt{(p-1)(1-\rho)}} (1, \dots, 1, \underset{i}{-(c-1)}, 0, \dots, 0)$$

$$e_p' = \frac{1}{\sqrt{(p-1)\rho}} (1, \dots, 1, \dots, 1, -(p-1))$$

- If  $\rho > 0$ , then

$$Y_1 = e_1' \cancel{\mathbf{Z}} = \frac{1}{\sqrt{P}} \sum_{i=1}^P Z_i$$

$$\frac{\lambda_1^\rho}{\sum \lambda_i} = \frac{1 + (P-1)\rho}{P} \approx \rho \quad \text{for large } P \text{ or } \rho \text{ close to 1}$$

$(P \rightarrow \infty \text{ or } \rho \rightarrow 1)$

d. Example.

<1> Let  $X_1$  = reading speed

$X_2$  = reading comprehension

$X_3$  = arithmetic speed

$X_4$  = — comprehension

representing the scores of a seventh-grade child.

- Let the correlation matrix  $P$  be

$$P = \begin{bmatrix} 1 & 0.698 & 0.264 & 0.081 \\ 0.698 & 1 & -0.061 & 0.092 \\ 0.264 & -0.061 & 1 & 0.594 \\ 0.081 & 0.092 & 0.594 & 1 \end{bmatrix}$$

SAS code ~~8~~ 8-1

# Lecture #18

## 8.3 Summarizing Sample variation by Principal Components

### (a) Sample principal Components

(1) Let  $X_1 = \begin{pmatrix} X_{11} \\ \vdots \\ X_{1P} \end{pmatrix}, \dots, X_n = \begin{pmatrix} X_{n1} \\ \vdots \\ X_{nP} \end{pmatrix}$  be  $n$  independent observations from  $p$ -dimensional space.

Data :  $X = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$

population  $X$  with  $\bar{E}(X) = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}$ ,  $\text{cov}(X) = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1P} \\ \vdots & \ddots & \vdots \\ \sigma_{P1} & \cdots & \sigma_{PP} \end{bmatrix}$

### (2) Sample P. Cs

- Let  $\bar{x}$ ,  $S$  and  $R$  be sample mean, sample covariance and correlation matrix
- ~~Sample~~ by a similar argument of obtaining the population P.C.s, the sample P.C.s are

$$\hat{Y}_i = \hat{e}_i' \bar{x}, \quad i=1, \dots, p \quad \text{with } S(\hat{Y}_i) = \hat{\lambda}_i$$

where  $(\hat{\lambda}_i, \hat{e}_i)$  being the eigen value-eigen vector pairs of  $S$ , with  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p$

- $\hat{Y}_1 = \hat{e}_1' \bar{x}$  maximizes  $S(a' \bar{x})$  subject to  $a'a=1$   
 $\Leftrightarrow \max_{a \neq 0} \frac{a'Sa}{a'a} = \hat{e}_1' S \hat{e}_1 = \lambda_1 = S(\hat{Y}_1)$

$\hat{Y}_2 = \hat{e}_2' \bar{x}$  maximizes  $S(a' \bar{x})$  with  $\text{cov}(a' \bar{x}, \hat{e}_1' \bar{x}) = 0$   
and  $a'a=1$ . ( $\text{cov}(a' \bar{x}, \hat{e}_1' \bar{x}) = a'S \hat{e}_1 = \lambda_1 a'e_1 = 0$ )  
 $\Leftrightarrow \max_{a \perp e_1} \frac{a'Sa}{a'a} = \hat{e}_2' S \hat{e}_2 = \lambda_2 = S(\hat{Y}_2)$

$\hat{Y}_i = \hat{e}_i' X$  satisfies

$$\max_{\substack{\text{a.s.t.e.} \\ \text{a.e.}}} \frac{a'sa}{a'a} = \hat{e}_i' S \hat{e}_i = \lambda_i = S(\hat{Y}_i)$$

- $\hat{Y}_i$ 's are uncorrelated linear combination with (S.P.C.)

largest sample variance.

$$\text{cov}(\hat{Y}_i, \hat{Y}_j) = 0, i \neq j \quad V(\hat{Y}_i) = \hat{\lambda}_i$$

$$\text{Total sample variance} = \sum_{i=1}^P S_{ii} = \sum_{i=1}^P \hat{\lambda}_i$$

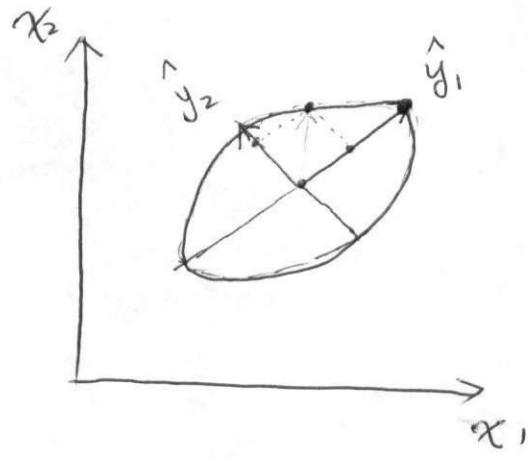
= Total sample variance due to sample principal components.

$\frac{\hat{\lambda}_i}{\sum \hat{\lambda}_i}$  = the proportion of total sample variance explained by the sample principal component.

$$\hat{r}_{y_i x_j} = \frac{SC(\hat{Y}_i, x_j)}{\sqrt{SC(\hat{Y}_i)} \sqrt{S(x_j)}} = \frac{\hat{e}_i' S \hat{x}_j}{\sqrt{\hat{\lambda}_i} \sqrt{S_{jj}}} = \frac{\lambda_i e_{ij}}{\sqrt{\lambda_i} \sqrt{S_{jj}}} = \frac{e_{ij} \sqrt{\lambda_i}}{\sqrt{S_{jj}}}$$

- Geometric interpretation of sample principal components

$$(X - \bar{X})' \sum_{i=1}^P \lambda_i e_i e_i' (X - \bar{X}) \\ = \sum_{i=1}^P \lambda_i (\hat{e}_i'(X - \bar{X}))^2 = C^2$$



① the ellipsoid  
on  $(X - \bar{X})' S^{-1}(X - \bar{X}) = C^2$

②  $e_i'(X - \bar{X})$  is on  
the direction of  
 $\hat{y}_i$

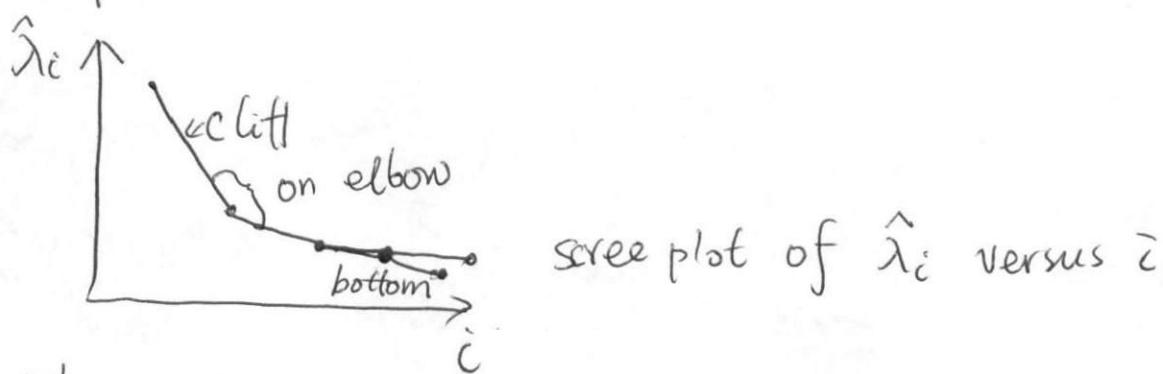
i.e. the direction of  
 $\hat{e}_i$

- Observation  $x_1, \dots, x_n$  are often "centered" by subtracting  $\bar{x}$

$$\hat{y}_i = \hat{e}_i' (x_1 - \bar{x}) = e_{i1}(x_{11} - \bar{x}_1) + \dots + e_{ip}(x_{p1} - \bar{x}_p)$$

(b) How many principal components are sufficient to explain sample variation?

- No definite answer. Need to look at
  - the amount of total sample variance explained
  - The relative sizes of eigen values.
  - subject-matter interpretations of the components
- Scree plot to help determine the number of principal components



The number of components is taken to be the point at which the remaining eigenvalues are relatively small and all about the same size such a point is usually near an elbow.

- Example 8.4

In study the size and shape relationship of male turtles,  
 The length, width and height of carapaces are measured.  
 The natural log of the dimensions of 24 male turtles have  
 been measured.

$$\hat{y}_i = \hat{e}'_i (X - \bar{X})$$

$$S = 10^{-3} \begin{bmatrix} 11.072 & 8.019 & 8.160 \\ 8.019 & 6.417 & 6.005 \\ 8.160 & 6.005 & 6.773 \end{bmatrix}$$

$$\bar{X}' = [4.725, 4.478, 3.703]$$

$$\hat{\lambda}_1' = 23.3 \times 10^{-3} \quad \hat{e}_1' = (0.683, 0.510, 0.523)$$

$$\hat{\lambda}_2' = 0.60 \times 10^{-3} \quad \hat{e}_2' = (-0.159, -0.594, 0.788)$$

$$\hat{\lambda}_3' = 0.36 \times 10^{-3} \quad \hat{e}_3' = (-0.713, 0.622, 0.324)$$

$$\frac{\hat{\lambda}_1}{\sum \hat{\lambda}_i} = 0.9605 \quad - \quad \frac{\hat{\lambda}_2}{\sum \hat{\lambda}_i} = 0.0247 \quad \frac{\hat{\lambda}_3}{\sum \hat{\lambda}_i} = 0.0148$$

$$\hat{y}_1 = \ln(\text{length}^{0.683} \cdot \text{width}^{0.51} \cdot \text{height}^{0.523}) - 7.447624$$

measuring the ln (volume) of carapace (with adjusted dimensions)

SAS code 8-2.txt

(c) Sample P.C. from standardized observations

1. Variables measured on different scales or with widely different ranges are often standardized.

- Let  $\tilde{X}_j = \begin{pmatrix} \frac{x_{j1} - \bar{x}_1}{\sqrt{s_{11}}} \\ \vdots \\ \frac{x_{jp} - \bar{x}_p}{\sqrt{s_{pp}}} \end{pmatrix} = D^{-\frac{1}{2}}(X_j - \bar{X})$   $D = \text{diag}(S) = \begin{pmatrix} s_{11} & & \\ & \ddots & \\ & & s_{pp} \end{pmatrix}$

$$Z = \begin{pmatrix} z_1' \\ \vdots \\ z_n' \end{pmatrix} = D^{-\frac{1}{2}} \begin{pmatrix} (x_1 - \bar{x})' \\ \vdots \\ (x_n - \bar{x})' \end{pmatrix} = D^{-\frac{1}{2}}(X - \bar{X})'$$

$$S(Z) = D^{-\frac{1}{2}} S D^{-\frac{1}{2}} = R(X)$$

2. The sample P.C.s of the standardized data  $Z$

are

$$\hat{Y}_i = \hat{e}_i' Z \quad i=1, \dots, p, \quad S(\hat{Y}_i) = \hat{\lambda}_i$$

where  $(\hat{\lambda}_i, \hat{e}_i)$  are the eigen value - vector pairs

of  $R$ ,  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$

- $S(\hat{Y}_i, \hat{Y}_j) = 0, \quad i \neq j \quad S(\hat{Y}_i) = \hat{\lambda}_i$

- $\sum_{i=1}^p \hat{\lambda}_i = P = \sum_{i=1}^p Y_{ii}$

- $\hat{y}_{ik}, z_k = e_{ik} \sqrt{\lambda_i} \quad i=1, \dots, p, \quad k=1, \dots, n$

(3) proportion of sample variance due to principal components

- $i^{\text{th}}$ :  $\frac{\hat{\lambda}_i}{P}$

- Scree plot

Select  $\hat{\lambda}_i$  with  $\hat{\lambda}_i > 1$

- extremely small  $\hat{\lambda}_p$ , however, may indicate linear dependency in variables

(e.g.:  $X_p = X_1 + \dots + X_{p-1}$  delete  $X_p$  before analysis)

## 8.4 Graphing the principal component

a) objective

1. check normality assumption

## code8-2

```
options ls=80 ps=47 nodate nonumber;
title1 h=2 'Male turtle shape data';
title2 h=1 'JW: E8-4';

data turtle;
infile 'Z:\\My Documents\\Teaching\\Stat524\\Fall 2010\\Data set\\E8-4.dat' firstobs=1;
input length width height;
lglength=log(length);
lgwidth=log(width);
lgheight=log(height);
run;

proc print data=turtle;
run;

title2 "Plot of raw data by length and height";
%plotit(data=track,labelvar=country,plotvars=m100 m200, color=black, colors=blue);
%plotit(data=turtle, plotvars=lglength lgheight);
run;

*ods html file="sas83.html";
ods graphics on;

proc princomp data=turtle cov out=pcturtle;
var lglength lgwidth lgheight;
run;

ods graphics off;
*ods html close;

proc print data=pcturtle;
run;

title2 "Scatter plot of the first 2 principal components";
%plotit(data=ptrack, labelvar=country, plotvars=prin2 prin1, color=black,
colors=blue);

%plotit(data=pcturtle, plotvars=prin2 prin1);
```

# Lecture #19

## 8.4 Graphing the principal components

### a) Objectives

- <1> check normality assumption
- <2> Identify suspect observations

### b) Checking normality assumption

- <1> Q-Q plot

• If population  $X \sim N_p(\mu, \Sigma)$ , then the  $i$ th P.C.

$$Y_i = e_i' X = e_{i1}x_1 + \dots + e_{ip}x_p$$

is also normal

• For a given sample  $\bar{X}_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1p} \end{pmatrix}, \dots, \bar{X}_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{np} \end{pmatrix}$  from  $\bar{X}$

We have <For each  $i=1, \dots, p$ ,  $\hat{Y}_i = \begin{pmatrix} \hat{y}_{i1} \\ \vdots \\ \hat{y}_{in} \end{pmatrix}$  with

$$\hat{Y}_{ij} = \hat{e}_i' \bar{X}_j = \hat{e}_{i1}x_{j1} + \dots + \hat{e}_{ip}x_{jp} \quad j=1, \dots, n$$

which should also be normal.

• We can construct Q-Q plot for each  $\hat{Y}_i$  to check normality

### (2) Scatter plot

• If  $x_1, \dots, x_n \sim N_p(\mu, \Sigma)$ , then  $\check{Y}_i, \check{Y}_j$  should be jointly normal

then we can construct scatter plot to check the normality

ac) Identify suspect observations

<1> expression for  $X_j$  ( $x_1, \dots, x_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$ )

- Let  $\hat{e}_1, \dots, \hat{e}_p$  be the eigen vectors of sample covariance matrix  $S$ . Let  $\hat{P} = (\hat{e}_1, \dots, \hat{e}_p)$ , then

$$\hat{P} \hat{P}' = I_{n \times p}$$

$$X_{j1} \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} + X_{j2} \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} + \dots + X_{jp} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} X_{j1} \\ \vdots \\ X_{jp} \end{pmatrix} = X_j$$

$$= \hat{P} \hat{P}' X_j$$

$$= (\hat{e}_1 \dots \hat{e}_p) \begin{pmatrix} \hat{e}_1' \\ \vdots \\ \hat{e}_p' \end{pmatrix} X_j$$

$$= (\hat{e}_1 \dots \hat{e}_p) \begin{pmatrix} \hat{y}_{j1} \\ \vdots \\ \hat{y}_{jp} \end{pmatrix}$$

$$= \hat{y}_{j1} \hat{e}_1 + \dots + \hat{y}_{jp} \hat{e}_p$$

(2) Identify suspect obs.

- If  $\hat{y}_1, \dots, \hat{y}_q$  can fit the obs very well

then  $\|D\|^2 = \|X_j - (\hat{y}_{j1} \hat{e}_1 + \dots + \hat{y}_{jq} \hat{e}_q)\|_2^2 \quad q \leq p$

$$= \|\hat{y}_{j,q+1} \hat{e}_{q+1} + \dots + \hat{y}_{jp} \hat{e}_p\|_2^2$$

$$= \hat{y}_{j,q+1}^2 + \dots + \hat{y}_{jp}^2$$

should be small.

- Suspect observations will often have at least one large coordinate among  $\hat{y}_{j_{q+1}}, \dots, \hat{y}_{j_p}$ .
- Strategy: Construct scatter plots and Q-Q plots for last few principal components to identify suspect observation.

$$\hat{y}_p = \hat{e}_{p1}x_1 + \dots + \hat{e}_{pp}x_p = \hat{e}_p' x$$

$$\hat{y}_{p+1} = \hat{e}_{p+1}x_1 + \dots + \hat{e}_{p+p}x_p = \hat{e}_{p+1}' x$$

(d) examples ~~8.7/8.4~~ 8.18

SAS PCA of national  
track records of  
Women

## 8.5 Large sample inference

a) Large sample properties of  $\hat{\lambda}_i$  and  $\hat{e}_i$

$\Leftrightarrow$  Large sample results for  $\hat{\lambda}_i, \hat{e}_i$

- Let  $x_1, \dots, x_n \stackrel{iid}{\sim} N_p(\mu, \Sigma)$

Let  $\lambda_1 > \lambda_2 \dots > \lambda_p > 0$  be eigen values of  $\Sigma$

Let  $\lambda = (\lambda_1, \dots, \lambda_p)'$  and  $\hat{\lambda} = (\hat{\lambda}_1, \dots, \hat{\lambda}_p)'$

$$\Lambda = \text{diag}(\lambda) = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \\ \vdots & \vdots \\ 0 & \lambda_p \end{pmatrix}$$

Then

$$(i) \quad \sqrt{n}(\hat{\lambda}_i - \lambda_i) \xrightarrow{d} N_p(0, 2\lambda^2)$$

$$(ii) \quad \sqrt{n}(\hat{e}_i - e_i) \xrightarrow{d} N_p(0, E_{-i})$$

$$E_{-i} = \lambda_i \sum_{j \neq i} \frac{\lambda_j}{\lambda_j - \lambda_i} e_j e_j'$$

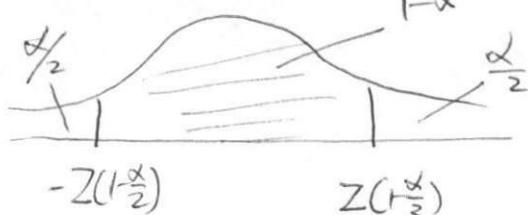
(iii)  $\hat{\lambda}_i$  and  $\hat{e}_{ij}$ ,  $j=1, \dots, p$  are independent,  $i=1, \dots, p$

(2) Confidence Interval for  $\lambda_i$

- By the above result, we have for large  $n$ ,

$$\hat{\lambda}_i \sim N(\lambda_i, \frac{2\lambda_i^2}{n})$$

Thus  $P\left(-\frac{|\hat{\lambda}_i - \lambda_i|}{\sqrt{\frac{2\lambda_i^2}{n}}} \leq Z(1-\frac{\alpha}{2})\right) = 1-\alpha$



Hence, 100(1- $\alpha$ )% C.I for  $\lambda_i$  is

$$-\frac{Z(1-\frac{\alpha}{2})}{\sqrt{\frac{1}{n}}} \leq \frac{\hat{\lambda}_i - \lambda_i}{\lambda_i \sqrt{\frac{2}{n}}} \leq Z(1-\frac{\alpha}{2})$$

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

(b) Testing for the equal correlation structure

(1) Large sample result

Let  $\Sigma_0 = \begin{bmatrix} \sigma_{11} & \dots & \sigma_{1P} \\ \vdots & \ddots & \vdots \\ \sigma_{P1} & \dots & \sigma_{PP} \end{bmatrix}$  and  $\sigma_{ij} = \sqrt{\sigma_{ii}} \sqrt{\sigma_{jj}} \rho_{ij}$  with  $\rho_{ij} = \rho$  for  $i, j = 1, \dots, P$

$$\text{Thus } \rho_0 = \begin{bmatrix} 1 & \dots & \rho \\ \vdots & \ddots & \vdots \\ \rho & \dots & 1 \end{bmatrix}$$

Let  $H_0: \rho = \rho_0$  v.s.  $H_1: \rho \neq \rho_0$

Let  $R = \begin{bmatrix} r_{11} & \dots & r_{1P} \\ \vdots & \ddots & \vdots \\ r_{P1} & \dots & r_{PP} \end{bmatrix}$  be the sample correlation matrix

$$\bar{r}_{-i} = \frac{1}{P-1} \sum_{j \neq i} r_{ij}$$

$$\bar{r}_{-} = \frac{1}{P} \sum_{i=1}^P \bar{r}_{-i}$$

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

$$\text{Then } T = \frac{n-1}{(1 - \bar{r}_{-})^2} \left[ \sum_{i < j} (r_{ij} - \bar{r}_{-})^2 - \hat{\gamma} \sum_{i=1}^P (\bar{r}_{-i} - \bar{r}_{-})^2 \right]$$

$$\sim \chi^2 \left( \frac{(P+1)(P-2)}{2} \right)$$

## (2) testing

- By the large sample result of T

We reject  $H_0$  at level  $\alpha$  if

$$T > \chi^2_{\frac{1}{(P+D(P-2))}} (1-\alpha)$$

When n is very large.

## code8-3

```

options ls=85 ps=65; *nodate nonumber;
title "SAS PCA of national track records of women";
data track;
  *infile 'Z:\\My Documents\\Teaching\\Stat524\\Fall 2010\\Data set\\T1-9.dat' firstobs=1
;
  input country$ m100 m200 m400 m800 m1500 m3000 marathon ;
  cards;
ARG    11.57   22.94   52.50   2.05    4.25    9.19   150.32
AUS   11.12   22.23   48.63   1.98    4.02    8.63   143.51
AUT   11.15   22.70   50.62   1.94    4.05    8.78   154.35
BEL   11.14   22.48   51.45   1.97    4.08    8.82   143.05
BER   11.46   23.05   53.30   2.07    4.29    9.81   174.18
BRA   11.17   22.60   50.62   1.97    4.17    9.04   147.41
.
.
.
;
run;

proc print data=track;
run;

title2 "Plot of raw data by m100 and m200";
%plotit(data=track,labelvar=country,plotvars=m100 m200, color=black, colors=blue);
run;

ods graphics on;
proc princomp data=track out=pctrack;
  *var m100 m200 m400 m800 m1500 m3000 marathon;
run;
ods graphics off;

proc sort data=pctrack;
  by prin1;
run;

proc print;
  id country;
  var prin1 prin2 m100 m200 m400 m800 m1500 m3000 marathon;
  title2 "Rankings by the 1st PC: Overall performance";
run;

proc sort data=pctrack;
  by prin2;
run;

proc print;
  id country;
  var prin1 prin2 m100 m200 m400 m800 m1500 m3000 marathon;
  title2 "Rankings by the 2nd PC: Short- vs long- range performance";
run;

title2 "Scatter plot of the first 2 principal components";
%plotit(data=pctrack, labelvar=country, plotvars=prin2 prin1, color=black,
colors=blue);

title2 "Scatter plot of the first 2 principal components";
%plotit(data=pctrack, labelvar=country, plotvars=prin7 prin6, color=black,
colors=blue);

proc univariate data=pctrack normal plot;
var prin7 prin6;
run;

```

code8-3

-----  
Assignment #6

The day assigned: Tuesday, November 2, 2010  
The day due: Thursday, November 11, 2010

Read Text: Chapters 8 and 9

1. Assigned problems: 8.20, 8.21, 8.22
2. Suggested additional problems: 8.10, 8.11, 8.13