

Chapter 9 Factor Analysis Lecture #20

9.1 Introduction

(a) what is factor analysis?

(1) Grouped variables

Sometimes, variables from groups

Within each group variables are highly related.

Variables between groups, however, have relatively small correlations. Thus it is reasonable to represent each group by a single characteristic or factor, which is responsible for the observed correlation

(2) Factor analysis

- One of the most widely used multivariate techniques introduced by Charles Spearman (1904) started as a controversial difficult subject (due to the lack of powerful computing tools). It has emerged as one of the most fascinating and useful tools.

Its applicability to many diverse fields such as, biology, chemistry, ecology, economics, education, political science, psychology, and sociology, etc. has been demonstrated.

- The main concern is to identify the internal relationships between a set of random variables and to describe the covariance structure in terms of a few underlying unobservable random quantities called factors.
- Spearman (1904) noticed a certain systematic pattern in a correlation matrix ρ of scores obtained on scholastic exams (of 30 preparatory school boys)

⊕
look
at book

$$\rho = \begin{matrix} C \\ F \\ E \\ M \\ D \\ MU \end{matrix} \left[\begin{array}{cccccc} 1.00 & 0.83 & 0.78 & 0.70 & 0.66 & 0.63 \\ 0.83 & 1.00 & 0.67 & 0.67 & 0.65 & 0.57 \\ 0.78 & 0.67 & 1.00 & 0.64 & 0.54 & 0.51 \\ 0.70 & 0.67 & 0.64 & 1.00 & 0.45 & 0.51 \\ 0.66 & 0.65 & 0.54 & 0.45 & 1.00 & 0.40 \\ 0.63 & 0.57 & 0.51 & 0.51 & 0.40 & 1.00 \end{array} \right]$$

Ignoring the diagonal, any two rows are almost proportional

$$C \rightarrow M \quad \frac{0.83}{0.67} \asymp \frac{0.78}{0.64} \asymp \frac{0.66}{0.45} \asymp \frac{0.63}{0.51} \asymp 1.2$$

He argued ρ can be explained by a model

$$X_i = L_i F + \varepsilon_i, \quad i=1, \dots, 6$$

X_i - i th variable observation

F - random variable representing general factor.

ε_i - specific factor

$$\text{COV}(X_i, X_j)$$

$$= \text{COV}(L_i F + \varepsilon_i, L_j F + \varepsilon_j)$$

$$= \text{COV}(L_i F, L_j F) + \text{COV}(\varepsilon_i, \varepsilon_j)$$

$$= (\cancel{L_i \text{COV}(F, F) L_j}) - L_i L_j V(F)$$

$$\rho(X_i, X_j) = \frac{L_i L_j V(F)}{\sqrt{L_i^2 V(F) + \sigma_i^2} \sqrt{L_j^2 V(F) + \sigma_j^2}}$$

(b) PCA versus FA

(1) PCA:

- An orthogonal transformation of the coordinate axes
(through the natural shape of the scatter plot of observations)
that partitions the total variance of all response into successively smaller partitions
- If the first few axes account for most of the total variance, and if their positions could be interpreted meaningfully, then the system could be described more parsimoniously by them
- Shortcomings:
 - 1) not invariant under transformation
 - 2) no rational criteria for selecting the number of PCs to retain

(2) FA:

- construct a fundamental model of covariance structure using some unobservable variable.
- Avoid the deficiency of PCA

9.2 The orthogonal factor model

a) orthogonal factor model with m common factors

1. the model

Let X be $p \times 1$ random vector with

$$E(X) = \mu \text{ and } \text{cov}(X) = \Sigma$$

Assume that the interrelationships between the elements of X can be explained by

$$X_{px1} = \mu_{px1} + L F_{p \times m \times 1} + \varepsilon_{px1}$$

Where $F = (F_1, \dots, F_m)'$ are called common factors,

which are random

$L = \{l_{ij}\}_{p \times m}$ are called factor loading matrix

l_{ij} : loading of i th variable on the j th factor

$\varepsilon = (\varepsilon_1, \dots, \varepsilon_p)'$ are called specific factors and random

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

⋮

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

2. Assumptions

We assume that

$$E(F) = 0, \quad \text{cov}(F) = I_m \quad E(\varepsilon) = 0, \quad \text{cov}(\varepsilon) = \Psi = \begin{pmatrix} \psi_{11} & 0 & \dots & 0 \\ 0 & \psi_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \psi_{pp} \end{pmatrix}$$

F and Σ are independent $\rightarrow \text{COV}(F, \Sigma) = 0$

(b) properties of and comments on the model.

1. population of the model

$$\bullet \quad \Sigma = \text{COV}(X) = \text{COV}(X - \mu) = E(LF + \Sigma)(LF + \Sigma)'$$

$$= E(LFF'L' + LF\Sigma' + \Sigma F'L' + \Sigma\Sigma')$$

$$= \underset{p \times m \text{ mp}}{LL'} + \Psi_{p \times p}$$

$$\bullet \quad \text{COV}(X, F) = \text{COV}(LF + \Sigma, F)$$

$$= L \text{COV}(F, F) + \text{COV}(\Sigma, F)$$

$$= L$$

$$\bullet \quad \sigma_{ii} = \text{var}(X_i)$$

$$= \text{var}(LF + \Sigma)_{ii}$$

$$= \cancel{\text{var}}(LL' + \Psi)_{ii}$$

$$= \sum_{j=1}^m l_{ij}^2 + \Psi_i$$

$$= h_i^2 + \Psi_i$$

\uparrow \nwarrow special variance
ith commonality

Statistically, $\Sigma = LL' + \Psi$ may not be solvable (when $m > p$)

(For $m=p$, $\Sigma = PAP'$)

$$\Sigma = \begin{bmatrix} 1 & 0.9 & 0.7 \\ 0.9 & 1 & 0.4 \\ 0.7 & 0.4 & 1 \end{bmatrix} \quad p=3$$

Let $m=1$,

$$= LL' + \Psi$$

$$= \begin{bmatrix} l_{11} \\ l_{21} \\ l_{31} \end{bmatrix} \begin{bmatrix} l_{11} & l_{21} & l_{31} \end{bmatrix} + \begin{bmatrix} \psi_1 & 0 & 0 \\ 0 & \psi_2 & 0 \\ 0 & 0 & \psi_3 \end{bmatrix}$$

$$= \begin{bmatrix} l_{11}^2 & l_{11}l_{21} & l_{11}l_{31} \\ l_{21}^2 & l_{21}l_{31} & \\ l_{31}^2 & & \end{bmatrix} + \begin{bmatrix} \psi_1 & 0 & 0 \\ 0 & \psi_2 & 0 \\ 0 & 0 & \psi_3 \end{bmatrix}$$

$$1 = l_{11}^2 + \psi_1 \quad 1 = l_{21}^2 + \psi_2$$

$$\begin{array}{ll} 0.9 = l_{11}l_{21} & 0.4 = l_{21}l_{31} \\ 0.7 = l_{11}l_{31} & 1 = l_{31}^2 + \psi_3 \\ \downarrow & \end{array}$$

$$\frac{l_{21}}{l_{11}} = \frac{4}{7}, \quad 0.9 = \frac{4}{7}l_{11}^2 \Rightarrow l_{11} = 1.575 \Rightarrow \psi_1 = 1 - (1.575)^2 = -0.575$$

$= V(\varepsilon_1)$
not much sense,

also $\text{COV}(X, F) = L$

$$l_{11} = \text{COV}(X_1, F)$$

$$= \text{corr}(X_1, F_1) \sqrt{\text{Var}(X_1)} \cdot \sqrt{\text{Var}(F)}$$

$$= \text{corr}(X_1, F_1) \Rightarrow l_{11} < 1 . \text{ again not make sense}$$

Solution may not be unique

$$\text{Let } TT' = I \text{ and } T'T = I$$

$$\Sigma = LL' + \Psi$$

$$= LTT'L' + \Psi$$

$$= (LT)(LT)' + \Psi$$

$$= L^* L^* + \Psi$$

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

$$= LTT'F + \Sigma$$

$$= L^* F^* + \Sigma \quad F^* = T'F \quad \text{with also } E(F^*) = 0$$

$$\text{and } V(F^*) = T'IT = T'T = I$$

Σ and F^* are also independent

9.3 Methods of Estimation

(a) The problem

1. Given $x_1, \dots, x_m \stackrel{iid}{\sim} X_{p \times 1}$, does the model

$$X = \mu + LF + \varepsilon, \quad L = (l_{ij})_{p \times m}, \quad F = (F_1, \dots, F_m)', \quad \text{var}(\varepsilon) = \begin{pmatrix} \Psi_1 & 0 \\ 0 & \Psi_p \end{pmatrix}$$

with $m < p$ adequately represent the data?

2. What are \hat{l}_{ij} , ψ_i and m ?

(b) The principal component method (PCM)

1. population case

$$\begin{aligned} \Sigma &= \sum_{i=1}^p \lambda_i e_i e_i' \\ &= (\sqrt{\lambda_1} e_1, \dots, \sqrt{\lambda_p} e_p) \begin{pmatrix} \sqrt{\lambda_1} e_1 \\ \vdots \\ \sqrt{\lambda_p} e_p \end{pmatrix} \end{aligned}$$

$$= LL'$$

$$\begin{cases} \hat{l}_{ij} = \sqrt{\lambda_j} l_{ij} \\ \psi_i = 0 \end{cases} \quad \text{is a solution}$$

$m = p$

but it is not practically useful. we hope $m \ll p$

$\Sigma = \lambda_1 e_1 e_1' + \dots + \lambda_m e_m e_m'$, if $\lambda_{m+1}^2 + \dots + \lambda_p^2$ is very small

$$L = [\sqrt{\lambda_1} e_1 \dots \sqrt{\lambda_m} e_m] \begin{bmatrix} \sqrt{\lambda_1} e_1 \\ \vdots \\ \sqrt{\lambda_m} e_m \end{bmatrix} \quad (m < p)$$

$$= LL'$$

$$\Sigma = LL' + \Psi \Rightarrow \begin{cases} \ell_{ij} = \sqrt{\lambda_j} e_{j|i} & i=1, \dots, p \\ \psi_i = \sigma_{ii} - \sum_{j=1}^m \ell_{ij}^2 & j=1, \dots, m \end{cases}$$

e_{ij} is the i th element of e_j

(2) Sample case

- $S = \sum_{i=1}^p \lambda_i \hat{e}_i \hat{e}_i'$ $S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$
- $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_p$

Pick $m < p$, then let

$$\tilde{L} = [\sqrt{\lambda_1} \hat{e}_1, \dots, \sqrt{\lambda_m} \hat{e}_m] = (\ell_{ij})_{p \times m} = (\sqrt{\lambda_j} \hat{e}_{j|i})_{p \times m}$$

$$\tilde{\Psi} = \begin{pmatrix} \tilde{\psi}_1 & 0 \\ 0 & \tilde{\psi}_p \end{pmatrix} \quad \text{with} \quad \tilde{\psi}_i = \tilde{s}_{ii} - \sum_{j=1}^m (\sqrt{\lambda_j} \hat{e}_{j|i})^2$$

$$= s_{ii} - \sum_{j=1}^m \hat{\lambda}_j \hat{e}_{j|i}^2$$

$$= s_{ii} - \sum_{j=1}^m \ell_{ij}^2$$

Lecture #21

(3) How to determine m ?

By prior knowledge

By proportion of total sample variance explained by factors F_1, \dots, F_n

$$\begin{aligned} 1) \quad S &= \sum_{i=1}^p \hat{\lambda}_i \hat{e}_i \hat{e}_i' \\ &= \sum_{i=1}^m \hat{\lambda}_i \hat{e}_i \hat{e}_i' + \sum_{i=m+1}^p \hat{\lambda}_i \hat{e}_i \hat{e}_i' \\ &\doteq \underline{L} \underline{L}' + \underline{\Psi} \\ &\quad \downarrow \\ &(\sqrt{\hat{\lambda}_1} \hat{e}_1, \dots, \sqrt{\hat{\lambda}_m} \hat{e}_m) \end{aligned}$$

2) Consider residual matrix

$$S - (\underline{L} \underline{L}' + \underline{\Psi})$$

The smaller the sum of squared entries of $S - (\underline{L} \underline{L}' + \underline{\Psi})$,

the better is the approximation

$$\Psi_{ij} = s_{ij} - \sum_{j=1}^m \hat{e}_{ij}^2$$

$$\begin{aligned} 3) \quad \text{Sum of squared entries of } S - (\underline{L} \underline{L}' + \underline{\Psi}) &= \sum_{i=m+1}^p \hat{\lambda}_i \hat{e}_i \hat{e}_i' = \hat{P}_2 \hat{\Lambda}_2 \hat{P}_2' \\ &\quad \cancel{\text{is}} \end{aligned}$$

$$\text{is almost } \text{tr} \left(\left(\sum_{i=m+1}^p \hat{\lambda}_i \hat{e}_i \hat{e}_i' \right) (\quad)' \right) \in \text{tr}((\quad)'(\quad))$$

$$= \text{tr}(\hat{P}_2 \hat{\Lambda}_2 \hat{P}_2')$$

$$= \text{tr}(\hat{\Lambda}_2^2) = \hat{\lambda}_{m+1}^2 + \dots + \hat{\lambda}_p^2$$

$$\left(\sqrt{\hat{\lambda}_{m+1}} \hat{e}_{m+1}, \dots, \sqrt{\hat{\lambda}_p} \hat{e}_{m+p} \right) (\quad)' (\quad) (\quad)$$

4^o) Thus if $\hat{\lambda}_{m+1}^2 + \dots + \hat{\lambda}_p^2$ is small, then the approximation is good, the model is adequate.

We normally Remark: The RMS of residual matrix ≤ 0.05 will be ideal case.

5^o Since

$s_{ii} = \sum_{j=1}^m \hat{e}_{ij}^2 + \hat{\psi}_i$, the contribution of the j th factor

F_j ($1 \leq j \leq m$) to s_{ii} is \hat{e}_{ij}^2

to $s_{11} + s_{22} + \dots + s_{pp} = \text{total sample variance}$ is

$$\hat{e}_{1j}^2 + \dots + \hat{e}_{pj}^2 = (\sqrt{\lambda_j} \hat{e}_{1j})^2 + \dots + (\sqrt{\lambda_j} \hat{e}_{pj})^2 \\ = \hat{\lambda}_j$$

6^o $\frac{\hat{\lambda}_j}{\hat{\lambda}_1 + \dots + \hat{\lambda}_p} = \frac{\hat{\lambda}_j}{s_{11} + \dots + s_{pp}}$ - proportion of total sample variance explained by the j th factor F_j .

$\frac{\hat{\lambda}_j}{p}$ if R is used.
 \uparrow
correlation matrix

$$m = \# \left(\hat{\lambda}_i > 1 \right)$$

rule of thumb.

SAS code 9-1.txt

Remark: $R = II' + \Psi \Leftrightarrow \Psi = 0$

(C) Principal factor method

↳ Population Case

- Assume that $\rho = LL' + \Psi$, then

1^o) ρ_{ij} is determined by $LL' (\hat{c} + \psi)$

$$2^o) \rho_{ij} = 1 + h_i^2 + \psi_i \quad (1 - \psi_i = h_i^2)$$

- Equivalently $\rho - \Psi = LL'$ (remove ψ_i from $\rho_{ii} = 1$)

(2) Sample case

- Suppose R is given and suppose that ψ_i can be estimated initially by ψ_i^*

Then we can replace $\rho_{ii} = 1$ by $1 + \psi_i^* = h_i^{*2}$ in R to obtain a "reduced" sample correlation matrix

$$R_r = \begin{bmatrix} h_1^{*2} & r_{12} & \cdots & r_{1P} \\ r_{21} & h_2^{*2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ r_{P1} & r_{P2} & \cdots & h_P^{*2} \end{bmatrix} \approx R - \tilde{\Psi}$$

- Factor R_r to get

$$R_r = \tilde{L}^* \tilde{L}^{*\top}$$

With $\tilde{L}_r^* = \{\tilde{e}_{ij}^*\}_{p \times m}$ if $m = P$, $R_r = \tilde{L}^* \tilde{L}^{*\top}$
0/w not equal

- Let $(\hat{\lambda}_i^*, \hat{e}_i^*)_{i=1, \dots, m}$ be the eigen pair of R_r , then the principal factor method

$$\mathbf{\hat{e}}^* = (\sqrt{\lambda_1^*} \hat{e}_1^*, \dots, \sqrt{\lambda_m^*} \hat{e}_m^*)$$

that is $\hat{e}_{ij}^* = \sqrt{\lambda_j^*} \cdot \hat{e}_{ij}^*$

$$\hat{\psi}_i^* = 1 - \sum_{j=1}^m \hat{e}_{ij}^{*2} = 1 - \hat{h}_i^{*2}$$

- Note that we start with initial estimates

$$h_i^{*2} \text{ (or } \hat{\psi}_i^* = 1 - h_i^{*2})$$

Then we get \hat{h}_i^{*2} (or $\hat{\psi}_i^*$)

Now we can start with \hat{h}_i^{*2} (or $\hat{\psi}_i^*$)

repeat the above step to obtain PF solutions

iteratively

(3) Prior values of h_i^{*2} (or $\hat{\psi}_i^*$) (estimates)

- Squared multiple correlation coefficient (SMC)

1° suppose $R = (r_{ij})$ $p \times p$. Then $r_{ii} = V_i F(X_i) = (1 - R^2)^{-1}$ (prior = SMC, in SAS)

$R_{\text{sample}}^2 = 1 - \frac{1}{r_{ii}}$ is called the sample multiple correlation coefficient.
Square of between x_i and $(p-1)$ remaining variables.

2° Let $\hat{\psi}_i^* = \frac{1}{r_{ii}}$ or $h_i^{*2} = 1 - \frac{1}{r_{ii}} = \text{SMC}_i$

- Input values for h_i^{*2} (or $\hat{\psi}_i^*$) (prior = INPUT)

$h_i^{*2} = 1 \rightarrow \text{PCM method}$

Selected random numbers for h_i^* (or ψ_i^*)
 (priors = random)

priors = max cm

$$h_i^* = \max_{|i \neq j|} |\gamma_{ij}|$$

(d) The maximum likelihood method.

<1> the method

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

$$\Sigma = LL' + \Psi$$

Then the MLE of L , Ψ and μ could be obtained
 by maximizing $L(\mu, \Sigma) = L(\mu, LL' + \Psi)$

$$= (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} e^{-\frac{1}{2} \text{tr} (\Sigma^{-1} (\sum_j (x_j - \bar{x}) (x_j - \bar{x})' + n(\bar{x} - \mu)(\bar{x} - \mu)'))}$$

subject to $L' \Psi^{-1} L$ being diagonal (uniqueness condition)

• (MLE method is scale invariant, we can also work with μ and ρ)

$$\Sigma = LL' + \Psi \quad V = \text{diag}(\Sigma)$$

$$\rho = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}} = V^{-\frac{1}{2}} (LL' + \Psi) V^{-\frac{1}{2}} = V^{-\frac{1}{2}} L (V^{\frac{1}{2}} L)' + V^{\frac{1}{2}} \Psi V^{\frac{1}{2}}$$

$L' \Psi^{-1} L$ being diagonal

$$(V^{\frac{1}{2}} L)' \cancel{V^{\frac{1}{2}}} \Psi^{-1} V^{\frac{1}{2}} \cancel{(V^{\frac{1}{2}} L)}$$

Hence,

Based on S , We have $\hat{\Sigma}$ and $\hat{\Phi}$ as MLE

Based on R (the standardized data),

We have $\hat{\Sigma}^{-\frac{1}{2}} \hat{\Lambda}$ and $\hat{\Sigma}^{-\frac{1}{2}} \hat{\Phi} \hat{\Sigma}^{-\frac{1}{2}}$ as MLE.

$$\hat{\Sigma}^{-\frac{1}{2}} = \text{diag}(\hat{\lambda}_1, \dots, \hat{\lambda}_p)$$

See proof in
supplement 9A

- The MLE of h_i^2 is $\hat{h}_i^2 = \sum_{j=1}^m \hat{\ell}_{ij}^2$ $1 \leq i \leq p$

- The proportion of total sample variance explained by the j th factor

$$\left\{ \begin{array}{l} \frac{\hat{\ell}_{1j}^2 + \dots + \hat{\ell}_{pj}^2}{S_{11} + \dots + S_{pp}} (= \lambda_j) \quad \text{for } S \\ \frac{\hat{\ell}_{1j}^2 + \dots + \hat{\ell}_{pj}^2}{P} (= \lambda_j^2) \quad \text{for } R \end{array} \right.$$

SAS 9-2. SAS

(2) Determination of the number of common factors

- A large sample test (likelihood ratio test)

1^o Bartlett's test

$$H_0: \Sigma = LL' + \Psi \quad H_1: \Sigma = \text{any other positive definite matrix}$$

By LR method,

$$\begin{aligned} L(\mu, \Sigma) &= \left(\frac{1}{2\pi}\right)^{\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2} \sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x}_j - \boldsymbol{\mu}) \right\} \\ &= (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2} \text{tr} \left[\Sigma^{-1} \left(\frac{n}{2} (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})' \right) \right] \right\} \end{aligned}$$

Under H_0 : $\hat{\boldsymbol{\mu}}_{MLE} = \bar{\mathbf{x}}$ and $\hat{\Sigma}_{MLE} = \hat{L}\hat{L}' + \hat{\Psi} = \hat{\Sigma}$

~~When~~ under $H_0 \cup H_1$, $\hat{\boldsymbol{\mu}}_{MLE} = \bar{\mathbf{x}}$, $\hat{\Sigma}_{MLE} = S_n = \frac{(n-1)S}{n}$

$$\begin{aligned} -2 \ln \Lambda &= -2 \ln \left[\left(\frac{|\hat{L}\hat{L}' + \hat{\Psi}|}{|S_n|} \right)^{-\frac{n}{2}} \cdot \frac{\exp \left\{ -\frac{n}{2} \text{tr} \left[(\hat{L}\hat{L}' + \hat{\Psi})^{-1} \cdot S_n \right] \right\}}{\exp \left\{ -\frac{1}{2} np \right\}} \right] \\ &= n \ln \left(\frac{|\hat{\Sigma}|}{|S_n|} \right) + n \left(\text{tr}(\hat{\Sigma}^{-1} S_n) - p \right) \quad \text{supplement 9A} \end{aligned}$$

$$= n \ln \left(\frac{|\hat{\Sigma}|}{|S_n|} \right)$$

\Leftarrow indicates $\text{tr}(\hat{\Sigma}^{-1} S_n) - p = 0$
if $\hat{\Sigma} = \hat{L}\hat{L}' + \hat{\Psi}$ is the MLE

Using Bartlett's correction,

$$(n-1-(2p+4m+5)/6) \ln \frac{|\hat{L}\hat{L}' + \hat{\Psi}|}{|S_n|} \Rightarrow \chi^2_{[(p-m)^2 - p-m]/2}$$

Hence, we reject H_0 if

$$(n-1-(2p+4m+5)/6) \ln \frac{|\hat{L}\hat{L}' + \hat{\Psi}|}{|S_n|} > \chi^2_{[(p-m)^2 - p-m]/2}$$

when both n and $n-p$ are large

$$(m < \frac{1}{2}(2p+1 - \sqrt{8p+1}))$$

Suggesting a higher-dimensional or nonlinear mechanism generating the correlation SAS (method = NL)

2°) Significant test results are common in testing the fit of a small number (m) ^(relative to n) of factors to correlation obtained from even moderately large sample.

~~Akaike's~~ Akaike Information Criterion (AIC)

1°) Alternatively, we may want to use Akaike's AIC instead of a sufficient test for determining the number of factors.

Under a m factors model : $\Sigma = LL' + \Psi$

$$\ln(L(m)) = C - \frac{n}{2} [\ln|LL' + \Psi| + \text{tr}(CL'L + \Psi)^{-1} S_n)]$$

$$\text{With } S_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

Then the AIC statistic is

$$\text{AIC}(m) = -2 \ln(L(m)) + [2p(m+1) - m(m-1)]$$

2°) the m -factor model with m corresponding to the smallest value of $\text{AIC}(m)$ is considered best

SAS 4.3.txt

code9-1

```
options ls=84 ps=55;
title 'SAS Factor Analysis example1';

data spearman (type=corr);
_type_='corr';
if _n_=1 then _type_='N';
infile cards missover;
input _name_$ c f e m d mu;
datalines;
N 33
c 1.0
f .83 1.0
e .78 .67 1.0
m .70 .67 .64 1.0
d .66 .65 .54 .45 1.0
mu .63 .57 .51 .51 .40 1.0
;

proc print data=spearman;
run;

ods html;
proc factor data=spearman method=prin res scree;
var c f e m d mu;
title2 'Princial Component Method';
run;

ods html close;

proc factor data=spearman method=prin res nfact=2;
title2 'PCM: Res matrix for nfact=2';
run;

proc factor data=spearman method=prin res nfact=3;
title2 'PCM: Res matrix for nfact=3';
run;

proc factor data=spearman method=prin res nfact=4;
title2 'PCM: Res matrix for nfact=4';
run;
```

code9-2

```
options ls=84 ps=55;
title 'SAS Factor Analysis Example 2';

data spearman (type=factor);
if _n_=1 then _type_='PRIORS';
else if _n_=2 then _type_='N';
else _type_='corr';
infile cards missover;
input _name_$ c f e m d mu;
datalines;
val .2 .3 .1 .1 .3 .2
N 33
c 1.0
f .83 1.0
e .78 .67 1.0
m .70 .67 .64 1.0
d .66 .65 .54 .45 1.0
mu .63 .57 .51 .51 .40 1.0
;

proc print data=spearman;
run;

proc factor data=spearman method=prin priors=input;
var c f e m d mu;
title2 'Princial Factor Method with input priors';
run;

proc factor data=spearman method=prin res nfact=2;
title2 'PCM: Res matrix for nfact=2';
run;

proc factor data=spearman method=prin res;* scree;
var c f e m d mu;
title2 'Princial Component Method';
run;

proc factor data=spearman method=prin res priors=smc;
title2 'PFM: priors=smc';
run;

proc factor data=spearman method=prin res priors=smc nfact=3;
title2 'PFM: priors=smc nfact=3';
run;

proc factor data=spearman method=prin res priors=max;
title2 'PFM: priors=max';
run;

proc factor data=spearman method=prin res priors=max nfact=3;
title2 'PFM: priors=max nfact=3';
run;
```

```

                                code9-3
options ls=80 ps=47 nodate nonumber;
title1 h=2 'SAS Factor Analysis Example 3';
title2 h=1 'Stock-Price Data (Weekly Rate of Return) JW: T8-4';

data stock;
*infile 'Z:\My Documents\Teaching\Stat524\Fall 2010\Data set\T8-4.dat' firstobs=1;
input JPMorgan Citibank WellsFargo Shell ExxonMobil;
datalines;
0.0130338      -0.0078431     -0.0031889     -0.0447693     0.0052151
0.0084862      0.0166886     -0.0062100     0.0119560     0.0134890
-0.0179153     -0.0086393     0.0100360     0.0000000     -0.0061428
0.0215589      -0.0034858     0.0174353     -0.0285917     -0.0069534
0.0108225      0.0037167     -0.0101345     0.0291900     0.0409751
0.0101713      -0.0121978     -0.0083768     0.0137083     0.0029895
0.0111288      0.0280044     0.0080721     0.0305433     0.0032290
.
.
.
0.0033740      -0.0153061     -0.0238245     -0.0016738     -0.0172270
0.0033626      0.0029016     -0.0030507     -0.0012193     -0.0097005
0.0170147      0.0095061     0.0181994     -0.0161758     -0.0075614
0.0103929      -0.0026612     0.0044290     -0.0024818     -0.0164502
-0.0127948     -0.0143678     -0.0187402     -0.0049759     -0.0163732
;
run;

proc print data=stock;
run;

proc factor method=prin res scree;
  title2 'Principal Component Method';
run;

proc factor method=prin n=2 priors=smc res;
  title2 'Principal Factor Method nfactor=2';
run;

/*Iterative method for estimation*/

proc factor method=prinIT priors=smc res;
  title2 'ITERATIVE Principal Factor Method';
run;

proc factor method=prinIT priors=smc n=2 res;
  title2 'ITERATIVE Principal Factor Method nfactor=2';
run;

proc factor method=ml res;
run;

proc factor data=stock method=ml n=1 res;
  title2 'MLE: nfactor=1';
run;

proc factor method=ml heywood n=2 res;
  title2 'MLE: nfactor=2';
run;

proc factor method=ml heywood nfact=3 res;
  title2 'MLE: nfactor=3';
run;

```

Lecture #22

9.4 Factor rotation

(a) Rotation

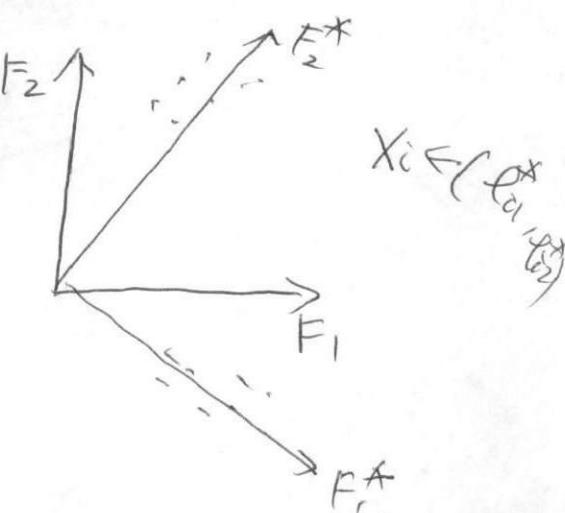
\Leftrightarrow nonuniqueness (might be a blessing)

- $S = \hat{\Sigma} \hat{\Sigma}' + \hat{\Psi}$
 $= \hat{\Sigma} T T' \hat{\Sigma}' + \hat{\Psi} \quad (TT' = I_m)$
 $= \hat{\Sigma}^* \hat{\Sigma}' + \hat{\Psi}$
- $\hat{\Sigma}$ and $\hat{\Sigma}^* = \hat{\Sigma} T$ represent the same covariance structure
 But it may be difficult to interpret using $\hat{\Sigma}$ and it may become easy to interpret factors using $\hat{\Sigma} T = \hat{\Sigma}^*$, a rotation of $\hat{\Sigma}$
- Ideally: After rotation, each variable loads highly on a single factor and has small or moderate loading on other factors

(2) groups variables

$$\begin{aligned} X - \mu &= LF + \varepsilon \\ &= LTT'F + \varepsilon \\ &= L^*F^* + \varepsilon \end{aligned}$$

If $m = 2$



(b) Select rotation (OR T)

<1> varimax rotation (orientation)

- Let $\hat{L}^* = \hat{L}T$ and let standardization ✓
 $\hat{e}_{ij}^{*2} = \frac{\hat{e}_{ij}^{*2}}{h_i^{*2}} = \frac{\hat{e}_{ij}^{*2}}{h_i^2}$ ← giving more weight to variables with small communality

Let $V = \sum_{j=1}^m \frac{1}{P} \left[\sum_{i=1}^P (\hat{e}_{ij}^{*2})^2 - \frac{\left(\sum_{i=1}^P \hat{e}_{ij}^{*2} \right)^2}{P} \right]$

$$\frac{1}{P} \sum_{i=1}^P \left[\hat{e}_{ij}^{*2} - \frac{1}{P} \left(\sum_{i=1}^P \hat{e}_{ij}^{*2} \right)^2 \right]^2$$

The variance of squares of scaled loading for j th factor
sample

- choose T such that V is maximized that is,
such that squares of scaled loadings spread out as much as possible.

- Equivalently, choose T such that among \hat{e}_{1j}^* , ..., \hat{e}_{Pj}^*
some are large, others are negligibly small.

(SAS: rotate = varimax)

(2) Quartimax rotation

- Let $\hat{L}^* = \hat{L}T$, choose T such that

$$\underbrace{\frac{1}{pm} \sum_{ij} \hat{e}_{ij}^{*4} - (\underbrace{\frac{1}{pm} \sum_{ij} \hat{e}_{ij}^{*2}}_{\text{Maximized}})^2}_{= \frac{1}{pm} \sum_i \hat{h}_i^2} \text{ is maximized}$$

(rotate = Quartimax)

(3) Non-orthogonal rotation (oblique rotation)

SAS (rotate = HK, rotate = promax)

9.5 Factor scores

(a) Estimation of factor scores

<1> The model

$$X - \mu = LF + \varepsilon \quad F = (F_1, \dots, F_m)'$$

↑ ↑ ↑
 Observable unobservable

$$S = \hat{L}\hat{L}' + \hat{\Psi}$$

(2) Factor scores

- Let f_j be the possible value F can take at j th observation ($j=1, \dots, n$)

The we estimate f_j by \hat{f}_j , which are called factor scores.

- Factor Scores can be useful in diagnostic processes and used as reduced data for any subsequent analysis

- To estimate f_j 's, we will treat $\hat{\ell}_j$ ($\hat{\ell}_{ij}^*$) and $\hat{\psi}_j$ as if they were the true values on the model.

(b) The weighted least square method

- Let $X - \mu = LF + \varepsilon$ $\text{COV}(\varepsilon) = \begin{pmatrix} \psi_1 & 0 \\ 0 & \psi_p \end{pmatrix} = \psi$

Assume that L , ψ and μ are known

$$\psi^{-\frac{1}{2}}(X - \mu - LF) = \psi^{-\frac{1}{2}}\varepsilon \quad \text{with } V(\psi^{-\frac{1}{2}}\varepsilon) = I_p$$

- choose f as the estimator of F such that

$$(\psi^{-\frac{1}{2}}\varepsilon)' (\psi^{-\frac{1}{2}}\varepsilon) = \sum_{i=1}^p \frac{\varepsilon_i^2}{\psi_i} = (X - \mu - LF)' \psi^{-1} (X - \mu - LF)$$

is minimized

(2) The solution (factor scores)

- Population case

$$\underbrace{\psi^{-\frac{1}{2}}(X - \mu)}_Y = \underbrace{\psi^{-\frac{1}{2}}LF + \psi^{-\frac{1}{2}}\varepsilon}_X$$

$$\hat{f} = (X'X)^{-1}X'Y$$

$$= (L'\psi^{-1}L)^{-1}L'\psi^{-1}(X - \mu)$$

- Sample case

1° Let \hat{L} and $\hat{\psi}$ and $\hat{\mu} = \bar{x}$ be the estimator of L , ψ and μ .

$$\hat{\Sigma} = \hat{L}\hat{L}' + \hat{\psi} = S$$

$$\hat{f}_j = (\hat{L}\hat{\psi}^{-1}\hat{L})^{-1}\hat{L}'\hat{\psi}^{-1}(x_j - \bar{x})$$

2° If we start with R then

$$\hat{f}_j = (\hat{\Sigma}_z' \hat{\Psi}_z^{-1} \hat{\Sigma}_z)^{-1} \hat{\Sigma}_z' \hat{\Psi}_z^{-1} \mathbf{z}_c^*$$

$$\text{Where } \mathbf{z}_c^* = (\text{diag}(S))^{-\frac{1}{2}} (X_c - \bar{X}) = D^{-\frac{1}{2}} (X_c - \bar{X})$$

$$\hat{\Sigma}_z = D^{\frac{1}{2}} \hat{\Sigma} D^{\frac{1}{2}}$$

- Comment

1° If $\hat{\Sigma}$ and $\hat{\Psi}$ are obtained via MLE, then

$\hat{\Sigma}' \hat{\Psi}^{-1} \hat{\Sigma}$ is a diagonal matrix (unique condition in MLE)

!!

$$\hat{\Sigma}_z' \hat{\Psi}_z^{-1} \hat{\Sigma}_z$$

2° If $\hat{\Sigma}$ and $\hat{\Psi}$ is obtained via PCM

and $\Psi_i = C$, $i=1, \dots, p$ then

\hat{f}_j are obtained essentially by unweighted (ordinary) least square method and

$$\hat{f}_j = (\hat{\Sigma}' \hat{\Sigma})^{-1} \hat{\Sigma}' (x_j - \bar{x})$$

$$\hat{\Sigma} = (\sqrt{\lambda_1} \hat{e}_1' \dots \sqrt{\lambda_m} \hat{e}_m')$$

$$= \begin{pmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_m} \end{pmatrix} \begin{pmatrix} \sqrt{\lambda_1} \hat{e}_1' \\ \vdots \\ \sqrt{\lambda_m} \hat{e}_m' \end{pmatrix} (x_j - \bar{x})$$

$$\hat{\Sigma}' = \begin{pmatrix} \lambda_1 & \dots & \lambda_m \end{pmatrix}$$

$$= \begin{pmatrix} \frac{1}{\sqrt{\lambda_1}} \hat{e}_1' (x_j - \bar{x}) \\ \vdots \\ \frac{1}{\sqrt{\lambda_m}} \hat{e}_m' (x_j - \bar{x}) \end{pmatrix}$$

← Scaled principal component

(C) The regression method

a) population:

$$\text{Assume } (F, \Sigma) \sim N_{m+p} (0, \begin{pmatrix} I_m & 0 \\ 0 & \Psi \end{pmatrix})$$

$$X - \mu = LF + \Sigma \sim N_{m+p} (0, LL' + \Psi)$$

$$(X - \mu, F)' \sim N_{m+p} (0, \begin{pmatrix} \Sigma = LL' + \Psi & L \\ L' & I_m \end{pmatrix})$$

By The result 4.6 $(X_1, X_2)' \sim N\left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}\right) \Rightarrow E(X_2 | X_1) = \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (X_1 - \mu_1)$

$$\begin{aligned} E(F | X - \mu) &= \mu_2 + \Sigma_{21} \Sigma_{11}^{-1} (X_1 - \mu_1) \\ x_2 &\quad x_1 \\ &= L' (LL' + \Psi)^{-1} (X - \mu) \end{aligned}$$

$$\text{cov}(F | X - \mu) = I - L' (LL' + \Psi)^{-1} L$$

$$\text{cov}(F | X - \mu) = I - L' (LL' + \Psi)^{-1} L$$

b) Sample solution $x_1, \dots, x_n \sim X \sim N_p (\mu, \Sigma) \quad \Sigma = LL' + \Psi$

$$\hat{f}_j^R = \hat{L}' (\hat{L} \hat{L}' + \hat{\Psi})^{-1} (x_j - \bar{x}) \quad (j=1, \dots, n) \quad \text{or} \quad \hat{f}_j^R = \hat{L}' S^{-1} (x_j - \bar{x})$$

If a correlation matrix is factored,

$$\hat{f}_j^R = \hat{L}_z' (\hat{L}_z \hat{L}_z' + \hat{\Psi}_z)^{-1} (z_j) \quad \text{or} \quad \hat{f}_j^R = \hat{L}_z' R^{-1} z_j$$

$$z_j = D^{-\frac{1}{2}} (x_j - \bar{x}) \quad D = \text{Diag}(S)$$

• Relationship between \hat{f}_j^R and \hat{f}_j^{LS}

$$\text{By } \hat{L}' (\hat{L} \hat{L}' + \hat{\Psi})^{-1} = (\hat{I} + \hat{L}' \hat{\Psi}^{-1} \hat{L})^{-1} \hat{L}' \hat{\Psi}^{-1} \quad \boxed{\text{See Exercise 9.6}}$$

$$\hat{f}_j^{LS} = (\hat{L}' \hat{\Psi}^{-1} \hat{L})^{-1} \hat{L}' \hat{\Psi} (x_j - \bar{x})$$

$$= (\hat{L}' \hat{\Psi}^{-1} \hat{L})^{-1} (\hat{I} + \hat{L}' \hat{\Psi}^{-1} \hat{L}) \hat{f}_j^R$$

$$= (\hat{I} + (\hat{L}' \hat{\Psi}^{-1} \hat{L})^{-1}) \hat{f}_j^R$$

If $\hat{L}, \hat{\Phi}$ are MLEs, $(\hat{L}' \hat{\Phi} \hat{L})^{-1} = \hat{\Delta}^{-1}$ where $\hat{\Delta}$ is a diagonal matrix. When $\hat{\Delta}$ is close to zero, \hat{f}_j^L is close to f_j^R

- If rotated loadings $\hat{L}^* = \hat{L}\hat{T}$ are used.

$$\hat{f}_j^* = T' \hat{f}_j, j=1, 2, \dots, n$$

(SAS. proc factor data = out =)

9.6 A strategy for factor analysis

(a) Most crucial decision in FA:

The number m of common factors

1. A large sample test (with method = ML)

- Suitable for approximately normal data points

- Rejection usually occurs for small n , large p .

(AIC) ~~(BIC)~~ criterion

2. the proportion of the total sample variance

explained by m factors

3. the RMS of residual matrix (≤ 0.05 , ideal case)

4. Subject-matter knowledge - prior knowledge

5. Interpretability

(b) Less crucial decision:

Estimation methods and rotations

1. Method = prin

= prin priors =

= m L

= image

= harris

= uls

= promax

= alpha

rotate = varimax

= quartimax

= equamax

= parsimax

= orthomax

= HK

= promax

= ~~oblimax~~

2. Successful FA usually leads to consistent results from different method.

(C) steps

1. PCFA

- plot factor scores check suspicious observations
(calculate the standardized scores and squared distance for each observation in assessing the normality)

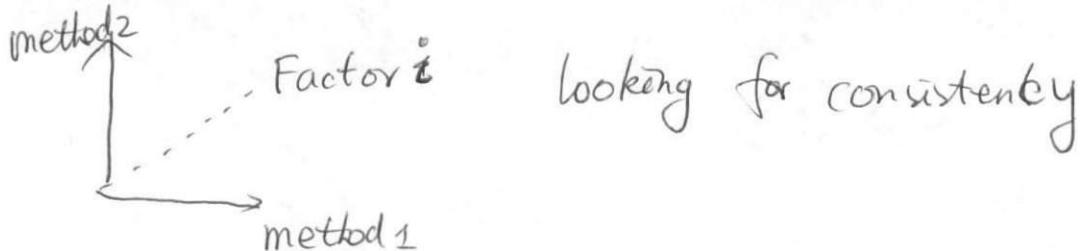
- try a varimax rotation

2. MLFA and varimax rotation

3. Perhaps other methods and rotations

④ Compare results obtained

4. • Compare factors (pattern of loadings)
- plot factor scores from different methods



5. ~~④~~ repeat above steps for different m to see if get better understanding and interpretability of the data.

6. For large data set (n large), split it into two and perform FA on each one, to check the consistency and stability of the result.

d) Example (Example 9.14)

```

        code9-4
options ls=80 ps=47 nodate nonumber;
title1 h=2 'SAS Factor Analysis Example 3';
title2 h=1 'Stock-Price Data (Weekly Rate of Return) JW: T8-4';

data stock;
*infile 'Z:\\My Documents\\Teaching\\stat524\\Fall 2010\\Data set\\T8-4.dat' firstobs=1;
input JPMorgan Citibank WellsFargo Shell ExxonMobil;
obs='*';
datalines;
0.0130338   -0.0078431   -0.0031889   -0.0447693   0.0052151
0.0084862   0.0166886   -0.0062100   0.0119560   0.0134890
-0.0179153   -0.0086393   0.0100360   0.0000000   -0.0061428
0.0215589   -0.0034858   0.0174353   -0.0285917   -0.0069534
0.0108225   0.0037167   -0.0101345   0.0291900   0.0409751
0.0101713   -0.0121978   -0.0083768   0.0137083   0.0029895
0.0111288   0.0280044   0.0080721   0.0305433   0.0032290
.
.
.
0.0033740   -0.0153061   -0.0238245   -0.0016738   -0.0172270
0.0033626   0.0029016   -0.0030507   -0.0012193   -0.0097005
0.0170147   0.0095061   0.0181994   -0.0161758   -0.0075614
0.0103929   -0.0026612   0.0044290   -0.0024818   -0.0164502
-0.0127948   -0.0143678   -0.0187402   -0.0049759   -0.0163732
;
run;

proc print data=stock;
run;

/*Factor Ration*/
proc factor data=stock method=ml heywood n=2 res rotate=varimax preplot plot;
  title2 'MLE: rotate=varimax';
run;

proc factor data=stock method=ml heywood n=2 res rotate=HK preplot plot;
  title2 'MLE: rotate=HK';
run;

/* Factor Scores*/
proc factor data=stock method=ml heywood n=2 res rotate=promax preplot plot
out=fscore;
  title2 'MLE: rotate=promax';
run;

proc print data=fscore;
run;

title2 "scatter plot of factor scores for factor1 and factor2";
%plotit(data=fscore,labelvar=obs, plotvars=factor2 factor1, color=black,
colors=blue);
run;

proc factor data=stock method=prin priors=asmc nfact=3 rotate=varimax out=fscore1;
run;

proc print data=fscore1;
run;

title2 "scatter plot of factor scores for factor1 and factor3";

```

```
code9-4
```

```
%plotit(data=fscore1,labelvar=obs, plotvars=factor3 factor1 ,color=black,
colors=blue);
run;

data fscorenew;
set fscore;
factor1mle=factor1;
run;

data fscorecompare;
merge fscorenew fscore1;
run;

%plotit(data=fscorecompare,labelvar=obs, plotvars=factor1mle factor1 ,color=black,
colors=blue);
run;
```

HW#7 9.19, 9.32

due Tuesday of Nov. 23, 2010

Lecture 23

chapter 11 Discrimination and classification

11.1 Objectives of discrimination and classification

(a) Discrimination (separation)

1. To set up a method of deciding whether the observations fall into groups and if so to delineate the groups

2. Examples

- IUPUI admission office
- banking and commercial finance
loan application

(b) classification (Allocation)

1. To construct a method of assigning a new observations to the correct groups (population)

(c) Connection of discrimination and classification

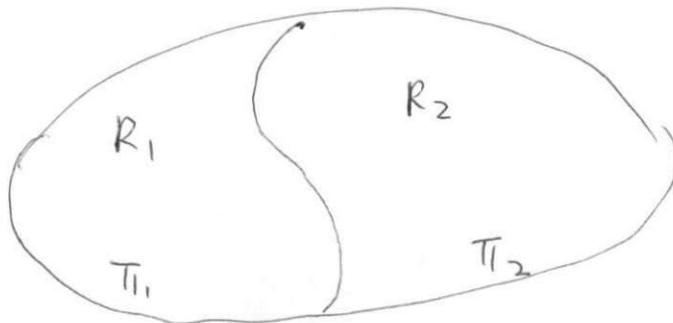
The two are interrelated, so closely, indeed that they are often confused

11.2 classification for the population

a) steps

1. rules are developed from "training" and "learning" sample
2. divide all possible sample outcomes into two groups

region → R_1 and R_2 , which correspond to populations π_1 and π_2 , respectively.



(3) There may not exist clear cut-off boundary
(no-error free allocation methods)

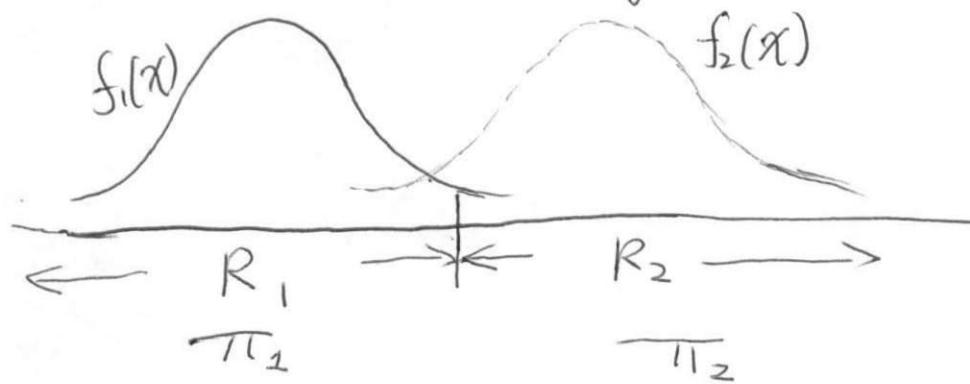
(b) goals

<1> to keep misclassification rate as small as possible
by taking the following into account.

(i) occurrence frequency from individual populations

(ii) misclassification cost (penalty)

(c) misclassification probability and cost



(1) Conditional Probability

- Conditional misclassification Prob

1° $P(\text{allocate } X \text{ to } \pi_1 \mid X \in \pi_2)$

$$= P(X \in R_1 \mid X \in \pi_2) = P(X \in R_1 \mid \pi_2) = P(1/2)$$

$$= \int_{R_1} f_2(x) dx$$

$\text{P}^{\text{2}} \text{ P}(\text{allocate } X \text{ to } \pi_2 \mid X \text{ is from } \pi_1)$

$$= P(2|1) = \int_{R_2} f_i(x) dx$$

• misclassification probability

1^o) $P(X \text{ is misclassified as } \pi_1)$

= $P(X \text{ is from } \pi_2, \text{ but allocated to } \pi_1)$

$$= P(X \text{ is from } \pi_2) \cdot P(1|2)$$

$$= P_2 \cdot P(1|2)$$

2^o) $P(X \text{ is misclassified as } \pi_2)$

$$= P_1 P(2|1) \quad (P_1 = 1 - P_2)$$

3^o) $P(X \text{ is correctly allocated to } \pi_1)$

$$= P_1 P(1|1)$$

$P(X \text{ is correctly allocated to } \pi_2)$

$$= P_2 P(2|2)$$

(2) Cost of misclassification
classified as

		π_1	π_2
		0	$C(2 1)$
true	π_1	0	
	π_2	$C(1 2)$	0

(d) classification rules

<1> Minimize expected (average) cost of misclassification
(MECM)

- Expected (average) cost of misclassification

$$\begin{aligned}
 ECM &= C(2|1) P_1 R(z|1) + C(1|2) P_2 R(1|2) \\
 &= C(2|1) P_1 \int_{R_2} f_1(x) dx + C(1|2) P_2 \int_{R_1} f_2(x) dx \\
 &= C(2|1) P_1 \int_{R-R_1} f_1(x) dx + C(1|2) P_2 \int_{R_1} f_2(x) dx \\
 &= C(2|1) P_1 + \int_{R_1} (C(1|2) P_2 f_2(x) - C(2|1) P_1 f_1(x)) dx
 \end{aligned}$$

- Regions R_1 and R_2 that minimizes ECM

1°) General case

$$R_1 : C(1|2) P_2 f_2(x) - C(2|1) P_1 f_1(x) \leq 0$$

$$R_2 : C(1|2) P_2 f_2(x) - C(2|1) P_1 f_1(x) > 0$$

$$\Leftrightarrow R_1 = \left\{ x : \frac{f_1(x)}{f_2(x)} \geq \frac{C(1|2)}{C(2|1)} \cdot \frac{P_2}{P_1} \right\}$$

$$R_2 = \left\{ x : \frac{f_1(x)}{f_2(x)} < \frac{C(1|2)}{C(2|1)} \cdot \frac{P_2}{P_1} \right\}$$

2°) Special cases

If $\frac{P_2}{P_1}$ unknown, then assume $\frac{P_2}{P_1} = 1$

$$\Leftrightarrow R_1 : \frac{f_1(x)}{f_2(x)} \geq \frac{C(1|2)}{C(2|1)}$$

$$R_2 : \frac{f_1(x)}{f_2(x)} < \frac{C(1|2)}{C(2|1)}$$

If $\frac{C(1|2)}{C(2|1)}$ is ~~unknown~~ unknown, then assume $\frac{C(1|2)}{C(2|1)} = 1$

$$\Rightarrow R_1 : \frac{f_1(x)}{f_2(x)} \geq \frac{P_2}{P_1}$$

$$R_2 : \frac{f_1(x)}{f_2(x)} < \frac{P_2}{P_1}$$

If both $\frac{P_2}{P_1}$ and $\frac{C(1|2)}{C(2|1)}$ unknown, assume that

$$\frac{P_2}{P_1} \cdot \frac{C(1|2)}{C(2|1)} = 1 \Rightarrow \begin{cases} R_1 : & f_1(x) \geq f_2(x) \\ R_2 : & f_1(x) < f_2(x) \end{cases}$$

Question: what if $f_i(x)$ are unknown. $i=1, 2$

(2) Other classification rules:

- Minimize total prob of misclassification (MTPM)

$$TPM = P_1 P(2|1) + P_2 R(1|2)$$

$$= P_1 \int_{R_2} f_1(x) dx + P_2 \int_{R_1} f_2(x) dx$$

$$= P_1 + \int_{R_1} P_2 f_2(x) - P_1 f_1(x) dx$$

$$\Rightarrow \begin{cases} R_1 : & \frac{f_1(x)}{f_2(x)} \geq \frac{P_2}{P_1} \\ R_2 : & \frac{f_1(x)}{f_2(x)} < \frac{P_2}{P_1} \end{cases}$$

Code:

11.3 classification with two multivariate normal distribution

(a) Equal covariance matrices $\Sigma_1 = \Sigma_2 = \Sigma$

(1). classification regions R_1 and R_2 based on MECM

- Densities

$$1^{\circ} f_i(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{(x-\mu_i)' \Sigma^{-1} (x-\mu_i)}{2}} \quad i=1,2$$

$$\frac{f_1(x)}{f_2(x)} = e^{-\frac{\frac{1}{2}(x-\mu_1)' \Sigma^{-1} (x-\mu_1) + \frac{1}{2}(x-\mu_2)' \Sigma^{-1} (x-\mu_2)}{2}}$$

- Regions (general case)

$$1^{\circ} R_1: \frac{f_1(x)}{f_2(x)} \geq \frac{C(1/2) P_2}{C(2/1) P_1} \Leftrightarrow r$$

$$\Leftrightarrow Q \geq \ln r$$

$$\Leftrightarrow (\mu_1 - \mu_2)' \Sigma^{-1} X - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \geq \ln r$$

$$2^{\circ} R_2: \dots \dots \dots < \ln r$$

Linear (\bar{x}) classification rule

(2) classification in practice (we assume $\Sigma_1 = \Sigma_2 = \Sigma$, but Σ unknown)

- Estimation of μ_i and Σ

$$1^{\circ} \hat{\mu}_i = \bar{x}_i \quad i=1,2$$

$$\bar{x}_c = \frac{1}{n_c} \sum_{j=1}^{n_c} x_{cj}, \quad S_c = \frac{1}{n_c-1} \sum_{j=1}^{n_c} (x_{cj} - \bar{x}_c)(x_{cj} - \bar{x}_c)', \quad c=1,2$$

$$(2) \begin{aligned} x_{11}, \dots, x_{1n_1} &\sim N_p(\mu_1, \Sigma) \\ x_{21}, \dots, x_{2n_2} &\sim N_p(\mu_2, \Sigma) \end{aligned}$$

$$2^{\circ} \hat{\Sigma} = \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1+n_2-2} = S_{pooled}$$

- The rule

1^o Allocate x_0 to

$$\pi_1 \text{ if } (\bar{x}_1 - \bar{x}_2)' \hat{\Sigma}^{-1} x_0 - \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' \hat{\Sigma}^{-1} (\bar{x}_1 + \bar{x}_2) \geq \ln r$$

$$\pi_2 \text{ if } \dots \dots \dots < \ln r$$

2^o If $r = 1$ ($\frac{C(1|2) P_2}{C(2|1) P_1} = 1$) then

$$R_1: C(\hat{\Sigma}^{-1}(\bar{x}_1 - \bar{x}_2))' x_0 \geq \frac{1}{2} C(\hat{\Sigma}^{-1}(\bar{x}_1 - \bar{x}_2))' (\bar{x}_1 + \bar{x}_2)$$

$$\Leftrightarrow \hat{\alpha}' x_0 \geq \hat{\alpha}' \cdot \frac{\bar{x}_1 + \bar{x}_2}{2}$$

$$R_2: \hat{\alpha}' x_0 < \hat{\alpha}' \cdot \frac{\bar{x}_1 + \bar{x}_2}{2}$$

$$\begin{cases} \hat{\alpha} = \hat{\Sigma}^{-1}(\bar{x}_1 - \bar{x}_2) \\ = \text{Spotted } (\bar{x}_1 - \bar{x}_2) \end{cases}$$

discriminant
coefficients

If let $y = \hat{\alpha}' x$, then $y_0 = \hat{\alpha}' x_0$, $\bar{y}_1 = \hat{\alpha}' \bar{x}_1$, $\bar{y}_2 = \hat{\alpha}' \bar{x}_2$



$$y_{ij} = \hat{\alpha}' x_{ij}$$



Then

$$R_1: y_0 \geq \frac{1}{2} (\bar{y}_1 + \bar{y}_2)$$

$$R_2: y_0 < \frac{1}{2} (\bar{y}_1 + \bar{y}_2)$$

(B) Scaling

- Uniqueness problem

Any $c\hat{\alpha} = c \cdot \hat{\Sigma}^{-1}(\bar{x}_1 - \bar{x}_2)$, $c > 0$ can serve as $\hat{\alpha}$ in above rule

- Scale $\hat{\alpha}$. $\hat{\alpha}^* = \frac{\hat{\alpha}}{\|\hat{\alpha}\|} = \frac{\hat{\alpha}}{\sqrt{\hat{\alpha}' \hat{\alpha}}} \quad \text{recommended}$ (only if x variable have been standardised)

If this is not the case, great deal of care must be exerted in interpreting the result

Assess the relative importance $(\hat{\alpha}' x_i)$ of variables x_1, \dots, x_p as discriminators

compare to

code11-1

```
options ls=85 ps=65;
title1 'SAS DISCRIM example 1';
title2 h=1 'Salmon data JW: T11-2';
data salmon;
  infile 'Z:\\My Documents\\Teaching\\stat524\\Fall 2010\\Data set\\T11-2.dat' firstobs=1;
  input species gender fresh marine;
run;

proc print data=salmon;
run;

proc plot data=salmon;
  plot fresh*marine=species;
run;

title2 "scatter plot of fresh and marine";
%plotit(data=salmon,labelvar=_blank_, symvar=species,
plotvars=fresh marine, color=black, colors=blue);
run;

proc discrim data=salmon method=normal pool=test
  wcov pcov manova;
class species;
var fresh marine;
title2 'Testing the equality of normal population parameters';
run;
```