

### Factor Analysis

$$\underline{X} = \underline{\mu} + \underline{L}\underline{F} + \underline{\varepsilon}$$

↓

$$\text{Cov}(\underline{\varepsilon}) = \Psi$$

$$\text{Cov}(\underline{F}) = I$$

$$\text{Cov}(\underline{X}) = \underline{L}\underline{L}^T + \Psi$$

$$= \underline{\mu} + (\underline{L}\underline{Q})(\underline{Q}^T\underline{F}) + \underline{\varepsilon}$$

$$= \underline{\mu} + \underline{L}^* \underline{F}^* + \underline{\varepsilon}$$

Last time: Principal factors estimation.

- iterative algorithm computing  $\underline{L}$  based on  $\hat{\Psi}$  and updating  $\hat{\Psi}$  based on  $\hat{L}$ .

### Maximum Likelihood estimation.

R function: factanal

$$\text{Assume } \underline{X} \sim N_p(\underline{\mu}, \underline{L}\underline{L}^T + \Psi)$$

$$\text{MLE of } \underline{\mu} = \hat{\underline{\mu}} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i$$

$$\ln L(\underline{L}, \Psi) = -\frac{n}{2} [\ln(|\underline{L}\underline{L}^T + \Psi|) + \text{trace}((\underline{L}\underline{L}^T + \Psi)^{-1} \hat{\Sigma})] + \dots$$

loglikelihood

$$\text{Where } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^T = \frac{n-1}{n} S$$

- maximizing  $\ln L$  is not particularly easy.
  - need to assume  $\Psi$  has positive diagonal elements etc.
- likelihood function can be used to examine / test adequacy of r factor model.

Test  $H_0: C = \underline{L}\underline{L}^T + \Psi$  for some  $p \times r$   $\underline{L}$  and diagonal  $\Psi$ .

$H_1: C$  is arbitrary pos. def.

$$V = 2 \times \text{difference in max'd } \ln L$$

$$\sim \chi^2(V) \text{ where } V = \frac{1}{2}(p-r)^2 - \frac{1}{2}(p+r)$$

$$\text{e.g. } p=5, r=2, \text{ then } V = \frac{9}{2} - \frac{7}{2} = 1$$

~~$$p=5, r=3, V = \frac{4}{2} - \frac{8}{2} < 0$$~~

Observe  $V = V_{\text{obs}}$  then p-value  $\doteq P(\chi^2(V) \geq V_{\text{obs}})$

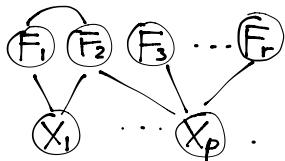
- small p-value indicates r factor model not adequate.
- larger p-value indicates r factor model is adequate.

### Choosing rotations

Graphical model representation

$$\underline{Y} = \begin{pmatrix} \underline{X} \\ \underline{E} \end{pmatrix} \quad \text{Cov}(\underline{Y}) = \begin{pmatrix} \underline{L}\underline{L}^T + \Psi & \underline{L} \\ \underline{L}^T & I \end{pmatrix} = C_+$$

Concentration matrix  $C_+ = \begin{pmatrix} \Psi^{-1} & -\Psi^{-1}L \\ -L^T\Psi^{-1} & I + L^T\Psi^{-1}L \end{pmatrix}$



Graphical structure is simplest if L has many 0 elements

$$L = \begin{pmatrix} 1_{11} & 1_{12} & \cdots & 1_{1r} \\ \vdots & & & \vdots \\ 1_{p1} & \cdots & 1_{pr} \end{pmatrix} = (\underline{l}_1 \cdots \underline{l}_r)$$

no link/edge between  $X_i$  &  $F_j$  if  $l_{ij}=0$ .

$\underline{l}_i^T \Psi^{-1} \underline{l}_j = 0$  then have no link between  $F_i$  &  $F_j$

$$\underline{l}_i = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \quad \underline{l}_k = \begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix} \quad \text{This occurs if } l_{ki} l_{kj} = 0 \text{ for all } k=1, \dots, p$$

### Type of rotations

① orthogonal  $L \rightarrow LQ$  where  $QQ^T = I$

② Oblique  $L \rightarrow LQ$  where Q is not orthogonal

Note: Assume variables standardized  $\rightarrow$  mean 0, variance 1  
 $\Rightarrow$  This allows loadings to be comparable.

### Variance rotation

Idea: Maximize variance of squared elements of each  $\underline{l}_1, \dots, \underline{l}_r$ .

maximize

$$\sum_{j=1}^r \sum_{i=1}^p \left[ \frac{l_{ij}^2}{\psi_i} - \frac{1}{P} \sum_{k=1}^p \frac{l_{kj}^2}{\psi_k} \right]^2$$

- forces squared elements of each  $\underline{l}_j$  to be as dispersed as possible.

- "encourages" 0 loadings

- varimax is the default in factanal.

### Oblique rotation

Take Q non-orthogonal with inverse  $Q^{-1}$ .  
Model:  $\tilde{X} = \tilde{\mu} + L\tilde{E} + \tilde{\xi} = \tilde{\mu} + (LQ)Q^{-1}\tilde{F} + \tilde{\xi}$   
 $= \tilde{\mu} + L^*F^* + \tilde{\xi}$

$$\text{Cov}(F^*) = Q^{-1} \underbrace{\text{Cov}(E)(Q^{-1})^T}_{I} = Q^{-1}(Q^{-1})^T \neq I$$

$F_1^*, \dots, F_r^*$  are no longer uncorrelated.

Problem? No practical reason to exclude correlated factors.

Example: Promax rotation.

- effectively raises varimax loadings to some power greater than 1 (e.g. 2, 4, 6)

$\Rightarrow$  emphasizes large loadings, reduces small loadings

### Application to athletics records data

p=8 events, n=55 countries.

Look at 1+2 factor models with various rotations (see Blackboard)

(r=1) uniqueness ( $\psi_1, \dots, \psi_8$ )

100m 200m 400m 800m 1500m 5000m 10000m marathon

$\psi_1, \dots, \psi_8 \rightarrow 0.447 \ 0.389 \ 0.242 \ 0.114 \ 0.059 \ 0.104 \ 0.367 \ 0.205$

loadings  $\rightarrow 0.743 \ 0.782 \ 0.871 \ 0.941 \ 0.970 \ 0.946 \ 0.794 \ 0.892$

But p-value =  $1.63 \times 10^{-16} \Rightarrow$  1 factor model is not adequate

(r=2)

100m 200m 400m 800m 1500m 5000m 10000m marathon

$\psi_1, \dots, \psi_8 \rightarrow 0.076 \ 0.123 \ 0.147 \ 0.128 \ 0.081 \ 0.037 \ 0.333 \ 0.092$

loadings  $\rightarrow 0.295 \ 0.378 \ 0.552 \ 0.711 \ 0.807 \ 0.904 \ 0.738 \ 0.913$

(varimax)  $0.915 \ 0.857 \ 0.740 \ 0.605 \ 0.517 \ 0.382 \ 0.350 \ 0.217$

p-value = 0.306  $\Rightarrow$  2 factor model is ok.

Blackboard: 2 other rotations

① rotation = "none"  $\rightarrow$  close to 1st 2 PC loadings

② promax.

### Cluster Analysis

Data  $x_1, \dots, x_n$  (often labelled by country or individual)

Goal: Try to group together (cluster) "similar" observations.

How?

Define distance function between pairs of points: Euclidean distance.  $(\sum_{i=1}^n (x_i - y_i)^2)^{\frac{1}{2}} = d(x, y)$

Manhattan distance  $d(x, y) = \sum_{i=1}^n |x_i - y_i|$

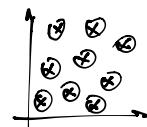
Describe inter-observation distances using a distance matrix.

$$D = \begin{pmatrix} 0 & d_{12} & d_{13} & \cdots & d_{1n} \\ d_{21} & \ddots & & & \vdots \\ \vdots & \ddots & \ddots & & \vdots \\ d_{n1} & \cdots & \ddots & 0 & \end{pmatrix} \text{ where } d_{ij} = d(x_i, x_j)$$

### Hierarchical clustering methods

Idea: Start with n clusters each with 1 observation.

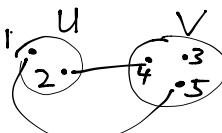
- Group together similar observations i.e. observations with small distances.



Need some way to measure distances between clusters

Three approaches:

- {
- ① Single linkage
- ② Complete linkage
- ③ Average linkage



① Single linkage:  $d(U, V) = d_{24}$   
 $d(U, V) = \min[d(x_i, x_j) : x_i \in U, x_j \in V]$

② complete linkage:  $d(U, V) = d_{15}$   
 $d(U, V) = \max[d(x_i, x_j) : x_i \in U, x_j \in V]$

③ average linkage:  $d(U, V) = \frac{1}{6}(d_{14} + d_{15} + d_{13} + d_{23} + d_{24} + d_{25})$   
 $d(U, V) = \text{average}[d(x_i, x_j) : x_i \in U, x_j \in V]$

Note:  $d(U, V) = d(V, U)$   
Triangle inequality?  $d(U, V) \leq d(U, W) + d(W, V)$

General algorithm:

1. Start with  $n$  clusters each with single observation  $\rightarrow$  distance matrix  $D$
2. Search  $D$  to find nearest pair of clusters (depends on approach), call these  $U+V$
3. Merge  $U+V$  into a single cluster and redefine  $D$ .
4. Repeat steps 2+3 until we have a single cluster

Example: Single linkage

$$D = \begin{pmatrix} 0 & & & & \\ 90 & 0 & & & \\ 3 & 370 & 0 & & \\ 4 & 6590 & 0 & 0 & \\ 5 & 110 & 280 & 0 & 0 \\ 1 & 2 & 3 & 4 & 5 \end{pmatrix}$$

Link observations 3 + 5 into a new cluster  
(call it "35")

$$\text{so } D' = \begin{pmatrix} 35 & 0 & & & \\ 11 & 0 & & & \\ 2 & 790 & 0 & & \\ 4 & 8650 & 0 & 0 & \end{pmatrix}$$

Link (35) & 1  $\Rightarrow$  (1, 3, 5)

$$D'' = \begin{pmatrix} 135 & 0 & & \\ 2 & 780 & 0 & \\ 4 & 0 & 0 & \end{pmatrix}$$

Link 2+4  $\rightarrow$  24

$$D''' = \begin{pmatrix} 135 & 0 & \\ 24 & 0 & \end{pmatrix}$$

This is called a Dendrogram:

