

FACTOR ANALYSIS - SOME CLASS NOTES

S. R. Searle

Biometrics Unit, Cornell University, Ithaca, New York

BU-748-M\*

November, 1981

Abstract

A description of factor analysis and its weaknesses.

---

Notation:  $\tilde{X}^*$  represents a vector of random variables. (The asterisk emphasizes that  $\tilde{X}^*$  is a vector, in contrast to the usual convention that  $\tilde{X}$  is a matrix.)

0. Introduction

Factor analysis is a method of attempting to explain the correlations that are found to exist among a set of random variables represented by the vector  $\tilde{X}^*$ . The attempted explanation is in terms of a number of hypothetical and unobservable random variables called factors. The underlying idea is that maybe the correlations existing among the variables we can observe can be explained by their being linear combinations of some fewer number of variables (that we cannot observe). Thus if  $\tilde{F}_{k \times 1}^*$  represents these unobservable random variables called factors, the population model for  $\tilde{X}^*$  is taken to be

$$\tilde{X}_{p \times 1}^* = \mu_{p \times 1} + L_{p \times k} \tilde{F}_{k \times 1}^* + \tilde{E}^*$$

---

\* In the Biometrics Unit Series, Cornell University, Ithaca, New York.

Before explaining this equation, the point must be made that principal components is not factor analysis - even though many people think it is. The concept of principal components never conjectures  $\tilde{X}^*$  as being explained in terms of any other variables - as does factor analysis. Principal components is concerned only with finding linear combinations of elements of  $\tilde{X}^*$  that explain maximum amounts of variation in the  $\tilde{X}^*$ 's, and never refers to a model involving unobservable variances - see Kendall and Lawley (1956).

The fiction that principal components and factor analysis are part of the same analysis, or worse, that principal components is just one way of doing factor analysis, comes about from the fact that both procedures do have in common the desire to explain the behaviour of observed variables. Furthermore, there is an approximation to the factor analysis model that permits of a solution using the principal components calculations. However, the lack of a model is no cause for not carrying out a principal components analysis, whereas the absolute need of a model for carrying out factor analysis is the feature that distinguishes the two analyses. Despite this, confusion still abounds as, for example, in Harman (1967), a widely used reference for factor analysis.

Also, many computer programs include both principal component and factor analysis in the same package - and this also engenders the false notion that they are the same thing. Indeed, prior to Joreskog's (1967) finding a workable algorithm for solving the maximum likelihood equations for factor analysis, this method of estimation was seldom used. What has been used, and probably still is being used, is what is sometimes called PFA, principal factor analysis. This is described in Section 4 of these notes. It is an approximation to factor analysis that makes use of the calculations of principal components analysis.

1. The Model

a. Description

The factor analysis model is

$$\tilde{X}_{p \times 1}^* = \underline{\mu} + \tilde{L}_{p \times k} \tilde{F}_{k \times 1}^* + \tilde{E}_{p \times 1}^* \quad (1)$$

where, using  $\mathcal{E}$  for expected value,

$$\underline{\mu} = \mathcal{E}(\tilde{X}^*) \quad , \quad (2)$$

and the other terms are defined as follows.

$\tilde{F}_{k \times 1}^*$  is a vector of  $k$  random variables,  $k < p$ , and these random variables are called factor scores or factors. They are unobservable.

$\tilde{L}_{p \times k} = \{l_{ir}\}$  is a  $p \times k$  matrix of coefficients to be determined that are called factor loadings.

$l_{ir}$  = loading of factor  $r$  in variable  $i$   
 = loading of variable  $i$  on factor  $r$  .

$\tilde{E}_{p \times 1}^*$  is a vector of  $p$  random (error) terms representing the difference between  $\tilde{X}^*$  and  $\underline{\mu} + \tilde{L}\tilde{F}^*$  .

Normality assumptions usually made are

$$\begin{bmatrix} \tilde{F}^* \\ \tilde{E}^* \end{bmatrix} \sim N \left[ \begin{pmatrix} \underline{0} \\ \underline{0} \end{pmatrix} \begin{pmatrix} \underline{U} & \underline{0} \\ \underline{0} & \underline{W} \end{pmatrix} \right] \quad (3)$$

so that

$$\tilde{X}^* \sim N(\underline{\mu}, \underline{V}), \text{ with } \underline{V} = \underline{L}\underline{U}\underline{L}' + \underline{W} \quad . \quad (4)$$

Simplified forms of  $\underline{U}$  and  $\underline{W}$  are often assumed:

$$\underline{U} = \underline{I} \quad (5)$$

and/or

$$\underline{W} = \underline{D}, \text{ diagonal} \quad (6)$$

or

$$\underline{W} = \sigma^2 \underline{I} \quad . \quad (7)$$

b. Rotations

Define

$$\underline{P}_{k \times k} \text{ as an orthogonal matrix .} \quad (8)$$

Then the model (1) is unchanged if written as

$$\underline{X}^* = \underline{\mu} + \underline{LP}'\underline{PF}^* + \underline{E}^* = \underline{\mu} + \underline{MG}^* + \underline{E}^* \quad (9)$$

for

$$\underline{M} = \underline{LP}' \quad \text{and} \quad \underline{G}^* = \underline{PF}^* \quad (10)$$

with

$$\begin{bmatrix} \underline{G}^* \\ \underline{E}^* \end{bmatrix} \sim N \left[ \begin{pmatrix} \underline{0} \\ \underline{0} \end{pmatrix}, \begin{pmatrix} \underline{PUP}' & \underline{0} \\ \underline{0} & \underline{W} \end{pmatrix} \right]$$

and

$$\underline{X}^* \sim N(\underline{0}, \underline{V})$$

with

$$\underline{V} = \underline{M}(\underline{PUP}')\underline{M}' + \underline{W} = \underline{LP}'\underline{PUP}'\underline{PL}' + \underline{W} = \underline{LUL}' + \underline{W} \text{ of (4) .}$$

Hence, in concept, the rotated model (9) is not different from the original model (1). Each is in terms of a set of factors,  $\underline{F}^*$  and  $\underline{G}^*$ , respectively, each uncorrelated with  $\underline{E}^*$  and in both cases  $\underline{X}^*$  has the same covariance structure, namely  $\underline{V} = \underline{LUL}' + \underline{W}$ . The only difference between the two models is that  $\underline{L}$  and  $\underline{M}$  do not have the same value. Since neither  $\underline{F}^*$ , nor  $\underline{G}^* = \underline{PF}^*$ , can be observed. This means that the factors represented by either  $\underline{F}^*$  or  $\underline{G}^*$  can be conceived of only to within the limits of an orthogonal transformation (or "rotation") - but this transformation is reflected in the coefficient matrices  $\underline{L}$  and  $\underline{M} = \underline{LP}'$ , which we seek to estimate.

Non-orthogonal transformations are also considered sometimes. These would involve considering the model (1) in the form

$$\tilde{X}^* = \tilde{\mu} + \tilde{L}\tilde{Q}^{-1}(\tilde{Q}\tilde{F}^*) + \tilde{E}^*$$

where  $\tilde{Q}$  is not orthogonal. Such transformations are sometimes called oblique rotations.

c. Numbers of factors

The whole objective of factor analysis is to explain the  $p$   $X^*$ -variables in terms of fewer  $F^*$ -variables, namely  $k$   $F^*$ -variables, with  $k < p$ . The variance structure given in (4) imposes limitations on  $k$ .

$$\tilde{V}_{p \times p} = \tilde{L}_{p \times k} \tilde{U}_{k \times k} \tilde{L}'_{k \times p} + \tilde{W}_{p \times p} \quad (11)$$

In truth we know none of the elements of this equation. But  $\tilde{X}^*$  can be observed and  $\tilde{V}$  estimated, so let us suppose we know  $\tilde{V}$ . Then, ignoring rank, (11) represents  $\frac{1}{2}p(p+1)$  equations in  $pk + \frac{1}{2}k(k+1) + \frac{1}{2}p(p+1) = \frac{1}{2}p(p+1) + \frac{1}{2}k(k+1+2p)$  unknowns, i.e., always more unknowns than equations. Clearly, even if  $\tilde{W}$  were known, there would be no way of solving (11). If  $\tilde{W}$  is confined to being a diagonal matrix than we have  $\frac{1}{2}p(p+1)$  equations in  $pk + \frac{1}{2}k(k+1) + p = \frac{1}{2}(k+1)(2p+k)$  unknowns, which permits of solution if

$$(k+1)(2p+k) \leq p(p+1)$$

i.e.,

$$k^2 + 2pk + k - p(p-1) \leq 0 \quad (12)$$

Since

$$k^2 + 2pk + k - p(p-1) = 0$$

when

$$k = \frac{1}{2}[-(2p+1) \pm \sqrt{(2p+1)^2 + 4p(p-1)}] = \frac{1}{2}[-(2p+1) \pm \sqrt{-(8p^2+1)}] ,$$

(12) is satisfied when

$$k < \frac{1}{2}[\sqrt{8p^2+1} - (2p+1)] ,$$

i.e.,

$$k < \sqrt{2p^2 + \frac{1}{4}} - (p + \frac{1}{2}) .$$

Examples

# of observable variables,  $p = 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 21 \ 100$   
 # of factors,  $k < 0 \ 0 \ 0 \ 2 \ 2 \ 2 \ 3 \ 3 \ 4 \ 4 \ 7 \ 41$  .

On taking  $\underline{U}$  as an identity matrix and  $\underline{W}$  as a diagonal matrix, as in (5) and (6), then (11) implies  $\frac{1}{2}p(p+1)$  equations in  $pk+p$  unknowns and we need

$$\frac{1}{2}(p+1) > k+1 ,$$

i.e.,

$$k < \frac{1}{2}(p-1) . \tag{13}$$

On further excluding orthogonal rotations (which for a  $k \times k$  matrix implies  $\frac{1}{2}k(k-1)$  relationships - orthogonality of  $k$  rows) we would then have  $\frac{1}{2}p(p+1)$  equations in  $pk - \frac{1}{2}k(k-1) + p$  unknowns requiring

$$\frac{1}{2}p(p+1) > pk - \frac{1}{2}k(k-1) + p ,$$

i.e.,

$$k^2 - 2pk - k + p^2 - p \geq 0 .$$

Since

$$k^2 - (2p+1)k + p(p-1) = 0$$

when

$$k = \frac{1}{2}[(2p+1) \pm \sqrt{(2p+1)^2 - 4p(p-1)}] = [p + \frac{1}{2} + \sqrt{2p + \frac{1}{4}}] ,$$

the inequality is satisfied when

$$k < p + \frac{1}{2} + \sqrt{2p + \frac{1}{4}} .$$

Examples

# of observed variables,  $p = 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \ 10 \ 20 \ 100$   
 # of factors,  $k < 0 \ 1 \ 2 \ 2 \ 3 \ 4 \ 4 \ 5 \ 6 \ 7 \ 14 \ 86$  .

Other than these results there seems to be inherently no part of factor analysis methodology that determines (or even gives guidance on) how many factors,  $k$ , there will be in a factor analysis representation of  $p$  observable random variables.

2. Maximum Likelihood Estimation

The factor analysis model is (1):

$$\tilde{X}_{p \times 1}^* = \underline{\mu} + \underline{L}\tilde{F}^* + \tilde{E}^* \quad (1)$$

$\tilde{X}^*$  represents p random variables that can be observed:  $\tilde{X}_{p \times N}$  represents a matrix of N realized values of  $\tilde{X}^*$ . Even though  $\tilde{F}^*$  cannot be observed, and despite there being no procedure for deciding how many elements are in  $\tilde{F}^*$  (other than some upper limit), the estimation problem is to estimate the factor loadings defined as L.

a. Equations to be solved

We confine attention to estimation by maximum likelihood and take assumptions (5) and (6) so that

$$\tilde{V} = \underline{L}\underline{L}' + \underline{D} \quad (15)$$

Since in the model (1) for  $\tilde{X}^*$  the random variables in  $\tilde{F}^*$  cannot be measured, the estimation problem is to estimate  $\underline{L}$ . And because  $\underline{L}$  occurs in  $\tilde{V}$ , as does  $\underline{D}$  also, the overall estimation problem is to estimate  $\underline{L}$  and  $\underline{D}$  of  $\tilde{V} = \underline{L}\underline{L}' + \underline{D}$ . The logical statistic for this purpose is the sum of squares and products matrix  $\tilde{S}$ ,

$$\tilde{NS} = \underline{X}\underline{X}' - N\underline{\bar{x}}\underline{\bar{x}}' \quad ,$$

where  $\underline{\bar{x}}$  is the vector of means  $\underline{\bar{x}} = \underline{X}\underline{1}/N$ . Then  $\tilde{S}$  has the Wishart distribution  $W(\tilde{V}, p, N)$  with density

$$\frac{|\tilde{S}|^{\frac{1}{2}(N-p-1)} \exp -\frac{1}{2}\text{tr}(\tilde{V}^{-1}\tilde{S})}{|\tilde{V}|^{\frac{1}{2}N} 2^{\frac{1}{2}NP} \prod_{i=1}^p \Gamma(\frac{n+1-i}{2})}$$

Therefore, for estimating elements of  $\underline{L}$  and  $\underline{D}$ , the log likelihood based on this density is used, namely

$$\theta = \text{constant} - \frac{1}{2}N \log |\tilde{V}| - \frac{1}{2}N \text{tr}(\tilde{V}^{-1}\tilde{S}) \quad .$$

This has to be minimized with respect to each  $d_r$  and each  $l_{ij}$ . The result (see Appendix) is

$$\sum_j \hat{l}_{ij}^2 + \hat{d}_i = s_{ii} \quad \text{for } i=1, \dots, p \quad (16)$$

and

$$\hat{\underline{L}}' \hat{\underline{D}}^{-1} \underline{S} = (\underline{I} + \hat{\underline{L}}' \hat{\underline{D}}^{-1} \hat{\underline{L}}) \hat{\underline{L}}' \quad (17)$$

with  $s_{ii}$  being the  $i$ 'th diagonal element of  $\underline{S}$ . Clearly, there is no analytic solution to (16) and (17).

b. No unique solution

Had we used  $\underline{M} = \underline{L}P'$  for  $P$  orthogonal these equations would be unchanged: (16), which for  $\underline{l}'_i$  being the  $i$ 'th row of  $\underline{L}$  and for  $\underline{m}'_i$  being the  $i$ 'th row of  $\underline{M}$ , would become

$$s_{ii} = \hat{m}'_i \hat{m}_i + \hat{d}_i = \hat{l}'_i P' P \hat{l}_i + \hat{d}_i = \hat{l}'_i \hat{l}_i + \hat{d}_i \quad ,$$

which is still (16); and (17) would be

$$\hat{\underline{M}}' \hat{\underline{D}}^{-1} \underline{S} = (\underline{I} + \hat{\underline{M}}' \hat{\underline{D}}^{-1} \hat{\underline{M}}) \hat{\underline{M}}' \quad ,$$

which is

$$\hat{\underline{P}} \hat{\underline{L}}' \hat{\underline{D}}^{-1} \underline{S} = \hat{\underline{P}} \hat{\underline{L}}' + \hat{\underline{P}} \hat{\underline{L}}' \hat{\underline{D}}^{-1} \hat{\underline{L}} \hat{\underline{P}} \hat{\underline{L}}'$$

and this reduces back to (17).

Thus we see that these maximum likelihood equations have no unique solution. They can be solved only to within an orthogonal transformation or rotation. Of this Lawley and Maxwell (p. 11, 1st ed.) write "In this situation all the statistician can do is to select a particular solution, one which is convenient to find, and leave the experimenter to apply whatever rotation he thinks desirable." Part of this advice ("a solution which is convenient to find") is familiar in the context of the general linear model not of full rank. But in that case the statistician can go on and tell the

experimenter how to usefully use his non-unique solution to yield unique results (estimates of estimable functions and so on). But this kind of advice is not possible in the factor analysis situation. Equations (16) and (17) have no unique solution, not because of any rank deficiency as in the linear model case, but simply by virtue of the nature of the model and analysis proposed. The model and analysis therefore seem open to question.

One way out of the predicament is to introduce more equations. Lawley and Maxwell do this by requiring  $\underline{L}$  to be such that

$$\underline{L}'\underline{D}^{-1}\underline{L} \text{ is diagonal} \quad . \quad (18)$$

However, justification for this seems to be slight and indeed they "... ignore the possibility that any of the diagonal elements of  $\underline{L}'\underline{D}^{-1}\underline{L}$  are equal", a fact which would negate the effectiveness of (18) ensuring a unique solution for (16) and (17).

### 3. Interpretation

#### a. The factor loadings

Having obtained  $\hat{\underline{L}}$ , a solution for (16) and (17), what can be done with it? The definition of  $\underline{L}$  given at the outset is

$$\underline{L} = \{\underline{l}_{ir}\}$$

for

$$\underline{l}_{ir} = \text{factor loading of variable } i \text{ on factor } r \quad .$$

In the words of Lawley and Maxwell, the experimenter, after obtaining one  $\hat{\underline{L}}$  is left "to apply whatever rotation he thinks desirable". This simply produces another  $\hat{\underline{L}}$  and, presumably, his ideas on what is "desirable" imply that one settles for an  $\hat{\underline{L}}$  that one likes the look of, meaning no doubt an  $\hat{\underline{L}}$  that one can interpret, and hopefully, interpret in a useful manner. Thus it is

that in books and papers on factor analysis one sees injunctions such as "rotate until all loadings, or as many as possible, are positive". And then, it seems, factors are interpreted by inspection. For example, suppose data on high school examinations in algebra, geometry, English and French are subjected to a factor analysis and the loadings come out as follows

	<u>Factor 1</u>	<u>Factor 2</u>	<u>Factor 3</u>
Algebra	.95	-.03	-1.01
Geometry	.87	.04	-1.04
English	-.01	.88	.96
French	.02	.87	.87

then it is not difficult to interpret the factors as (1) an aptitude in mathematics, (2) an aptitude in language, and (3) the difference between them. But this is based purely on the outcome of the data. Had factor 3 not been as shown, but had yielded loadings  $-.03$ ,  $.96$ ,  $-.46$  and  $.59$ , what kind of interpretation would that bring forth? And if one didn't like the factors so produced, or was unable to interpret them, is it then a statistically sound procedure to simply keep rotating until one does like the outcome and can interpret it? What assurance is there that future data will yield concomitant interpretation? Shouldn't modelling be done prior to data analysis and not post hoc?

Examples in the literature of conducting a factor analysis on "data" constructed from known factors, or with known relationships, are not reassuring: factor analysis does not always reveal the factors that are known to be the basis of the "data". Furthermore, the whole problem is additionally complicated by the question: how many factors? Allowing for differing numbers of factors does not always produce similar results: e.g., a four-factor analysis does not always yield factors in common with a three-factor analysis. As Press (p. 312) says, "... interpretation of results ... is an extremely important but often troublesome aspect of factor analysis." See also Armstrong (1967).

b. The factors

Interpretation of factor analysis is usually in terms of the relative magnitudes of the estimated factor loadings (elements of  $\hat{\underline{L}}$ ). Thus a study reported in Press (p. 314) consisted of a factor analysis of seven factors on the basis of 63 securities. Interpretation based on  $\hat{\underline{L}}$  of order  $63 \times 7$  directed attention initially to the fact that the entries in the first column of  $\hat{\underline{L}}$  were all large. The conclusion from this was that factor 1 affected each of the 63 variables to a large extent and could therefore be interpreted as some kind of overall market factor that affects all securities.

But why cannot the factors themselves be estimated directly? Lawley and Maxwell (p. 88, 1st ed.) suggest

$$\underline{\tilde{F}}^* = \underline{\hat{L}}^{-1} \underline{\tilde{S}}^{-1} \underline{\tilde{X}}^* .$$

However, since the model equation is

$$\underline{\tilde{X}}^* = \underline{\mu} + \underline{\tilde{L}} \underline{\tilde{F}}^* + \underline{\tilde{E}}^*$$

with  $\text{var}(\underline{\tilde{X}}^*) = \underline{\tilde{V}}$ , and because  $\hat{\underline{L}}$  will have full column rank (presumably  $\underline{\tilde{L}}$  does also, for otherwise  $\underline{\tilde{F}}^*$  would include variables that were utilized only as linear combinations of other variables), why not estimate  $\underline{\tilde{F}}^*$  as

$$\underline{\tilde{F}}^* = (\underline{\hat{L}}' \underline{\tilde{S}}^{-1} \underline{\hat{L}})^{-1} \underline{\hat{L}}' \underline{\tilde{S}}^{-1} (\underline{\tilde{X}}^* - \underline{\mu}) ,$$

analogous to GLS procedures; or perhaps as

$$\underline{\tilde{F}}^* = (\underline{\hat{L}}' \underline{\hat{D}}^{-1} \underline{\hat{L}})^{-1} \underline{\hat{L}}' \underline{\hat{D}}^{-1} (\underline{\tilde{X}}^* - \underline{\mu})$$

or, using LS ideas as

$$\underline{\tilde{F}}^* = (\underline{\hat{L}}' \underline{\hat{L}})^{-1} \underline{\hat{L}}' (\underline{\tilde{X}}^* - \underline{\mu})$$

or, using the ideas of predicting random variables from mixed models (see Searle, 1971, p. 462) as

$$\text{predicted } \tilde{F}^* = \tilde{L}'\tilde{V}^{-1}(\tilde{X}^* - \tilde{\mu}) \quad .$$

Any of these ways gives a matrix whose row elements indicate the extent to which the variables in  $\tilde{X}^*$  get utilized in forming the elements of  $\tilde{F}^*$ . In this way an indication is directly available as to the composition of  $\tilde{F}^*$ .

### c. Correlation

It is often argued that factor analysis should be free of the effects of scale and dispersion. The analysis is then carried out using the correlation matrix  $\tilde{R}$  in place of  $\tilde{S}$ . And to take account of intraclass correlations among the  $\tilde{X}$ 's, these are sometimes estimated as  $r_i$ . Then values  $1 - r_i$  for  $i = 1, \dots, p$  are used in place of the diagonal unities in  $\tilde{R}^2$ .  $r_i d_i$  is then called a communality.

### 4. Factor Analysis and Principal Components

Suppose that the error terms  $\tilde{E}^*$  in the factor analysis model (1) are ignored - or at least that their variances, the  $d_i$  of  $\tilde{D}$ , are deemed small enough to be ignored. Then in (15) the equation for  $\tilde{V}$  is

$$\tilde{V}_{p \times p} = \tilde{L}_{p \times k} (\tilde{L}')_{k \times p} \quad . \quad (19)$$

In principal components analysis we dealt with the eigen roots and vectors of  $\tilde{V}$  summarized in the equation

$$\tilde{B}'\tilde{V}\tilde{B} = \tilde{D}_\lambda$$

where  $\tilde{B}$  is an orthogonal matrix, its columns being eigen vector of  $\tilde{V}$ , and  $\tilde{D}_\lambda$  is a diagonal matrix of the corresponding eigen roots. Thus

$$\tilde{V} = \tilde{B}\tilde{D}_\lambda\tilde{B}' = \tilde{B}\tilde{D}_\lambda^{\frac{1}{2}}(\tilde{B}\tilde{D}_\lambda^{\frac{1}{2}})' \quad (20)$$

where  $(\tilde{D}_\lambda^{\frac{1}{2}})^2 = \tilde{D}_\lambda$ , which exists because  $\tilde{V}$  is p.d. and has positive eigen roots.

Comparing (19) and (20) we see that a possible value of  $\tilde{L}$  is

$$\tilde{L} = \tilde{B}\tilde{D}_\lambda^{\frac{1}{2}} \quad .$$



Appendix

Derivation of Maximum Likelihood Equations

We need to maximize

$$\theta = \text{constant} - \frac{1}{2}N \log|\underline{V}| - \frac{1}{2}N \text{tr}(\underline{V}^{-1}\underline{S}) \quad (\text{A1})$$

with respect to elements of  $\underline{L} = \{l_{ij}\}$  and of  $\underline{D} = \text{diag}\{d_r\}$  where

$$\underline{V} = \underline{L}\underline{L}' + \underline{D} \quad . \quad (\text{A2})$$

The following results from matrix differentiation are used:

$$\frac{\partial}{\partial \underline{x}} \log|\underline{A}| = \text{tr}\left(\underline{A}^{-1} \frac{\partial \underline{A}}{\partial \underline{x}}\right) \quad \text{and} \quad \frac{\partial}{\partial \underline{x}} \underline{A}^{-1} = -\underline{A}^{-1} \frac{\partial \underline{A}}{\partial \underline{x}} \underline{A}^{-1} \quad . \quad (\text{A3})$$

Also,

$$\frac{\partial \underline{A}}{\partial a_{ij}} = \underline{F}_{ij} = \left\{ \begin{array}{l} \text{null matrix except that} \\ \text{element (i,j) is unity} \end{array} \right\} \quad (\text{A4})$$

and

$$\text{tr}(\underline{A}\underline{F}_{ij}) = a_{ji} \quad . \quad (\text{A5})$$

First, from (A1),

$$\frac{\partial \theta}{\partial d_r} = \frac{-N}{2} \text{tr}\left(\underline{V}^{-1} \frac{\partial \underline{V}}{\partial d_r}\right) + \frac{N}{2} \text{tr}\left(\underline{V}^{-1} \frac{\partial \underline{V}}{\partial d_r} \underline{V}^{-1}\underline{S}\right)$$

and equating this to 0 gives

$$\begin{aligned} \text{tr}(\underline{V}^{-1}\underline{F}_{rr}) &= \text{tr}(\underline{V}^{-1}\underline{F}_{rr}\underline{V}^{-1}\underline{S}) \\ &= \text{tr}(\underline{V}^{-1}\underline{S}\underline{V}^{-1}\underline{F}_{rr}) \quad . \end{aligned}$$

Hence, using (A5)

$$\begin{aligned} (r,r)\text{'th element of } \underline{V}^{-1} &= (r,r)\text{'th element of } \underline{V}^{-1}\underline{S}\underline{V}^{-1} \\ \therefore (r,r)\text{'th element of } \underline{V} &= (r,r)\text{'th element of } \underline{S} \end{aligned}$$

$$\therefore \sum_j l_{rj}^2 + d_r = s_{rr} \quad . \quad (\text{A6})$$

Second, from (A1) again

$$\frac{\partial \theta}{\partial l_{ij}} = \frac{-N}{2} \operatorname{tr} \left( \tilde{V}^{-1} \frac{\partial \tilde{V}}{\partial l_{ij}} \right) + \frac{N}{2} \operatorname{tr} \left( \tilde{V}^{-1} \frac{\partial \tilde{V}}{\partial l_{ij}} \tilde{V}^{-1} \tilde{S} \right)$$

and equating this to 0 gives

$$\operatorname{tr} \left( \tilde{V}^{-1} \frac{\partial \tilde{L}\tilde{L}'}{\partial l_{ij}} \right) = \operatorname{tr} \left( \tilde{V}^{-1} \tilde{S} \tilde{V}^{-1} \frac{\partial \tilde{L}\tilde{L}'}{\partial l_{ij}} \right) . \quad (\text{A7})$$

Now, for any symmetric matrix  $\tilde{M}$

$$\begin{aligned} \operatorname{tr} \left( \tilde{M} \frac{\partial \tilde{L}\tilde{L}'}{\partial l_{ij}} \right) &= \operatorname{tr} (\tilde{M} \tilde{F}_{ij} \tilde{L}' + \tilde{M} \tilde{L} \tilde{F}_{ji}) \\ &= \operatorname{tr} (\tilde{L}' \tilde{M} \tilde{F}_{ij} + \tilde{M} \tilde{L} \tilde{F}_{ji}) \\ &= (j,i) \text{'th element of } \tilde{L}' \tilde{M} + (i,j) \text{'th element of } \tilde{M} \tilde{L} \\ &= 2[(i,j) \text{'th element of } \tilde{M} \tilde{L}] \quad . \end{aligned} \quad (\text{A8})$$

Therefore

$$\left\{ \operatorname{tr} \left( \tilde{M} \frac{\partial \tilde{L}\tilde{L}'}{\partial l_{ij}} \right) \right\} = 2 \tilde{M} \tilde{L} \quad .$$

Applying this to (A7) gives

$$\tilde{V}^{-1} \tilde{L} = \tilde{V}^{-1} \tilde{S} \tilde{V}^{-1} \tilde{L} \quad ,$$

i.e.,

$$\tilde{L} = \tilde{S} \tilde{V}^{-1} \tilde{L} \quad .$$

It is convenient to work this into an alternative form that does not involve  $\tilde{V}^{-1}$ .

$$\begin{aligned} \tilde{L}\tilde{L}' &= \tilde{S} \tilde{V}^{-1} \tilde{L}\tilde{L}' = \tilde{S} \tilde{V}^{-1} (\tilde{V} - \tilde{D}) = \tilde{S} - \tilde{S} \tilde{V}^{-1} \tilde{D} \\ \therefore \tilde{V} \tilde{D}^{-1} \tilde{L}\tilde{L}' &= \tilde{V} \tilde{D}^{-1} \tilde{S} - \tilde{S} \quad , \end{aligned}$$

and so

$$(\underline{\underline{L}}\underline{\underline{L}}' + \underline{\underline{D}})\underline{\underline{D}}^{-1}\underline{\underline{L}}\underline{\underline{L}}' = (\underline{\underline{L}}\underline{\underline{L}}'\underline{\underline{D}}^{-1} + \underline{\underline{I}})\underline{\underline{S}} - \underline{\underline{S}}$$

or

$$\underline{\underline{L}}\underline{\underline{L}}'\underline{\underline{D}}^{-1}\underline{\underline{L}}\underline{\underline{L}}' + \underline{\underline{L}}\underline{\underline{L}}' = \underline{\underline{L}}\underline{\underline{L}}'\underline{\underline{D}}^{-1}\underline{\underline{S}} \quad .$$

Then, using  $\underline{\underline{L}}\underline{\underline{L}}'P = \underline{\underline{L}}\underline{\underline{L}}'Q \Rightarrow \underline{\underline{L}}'P = \underline{\underline{L}}'Q$  for any real  $P, Q$  and  $\underline{\underline{L}}$ , we have

$$\underline{\underline{L}}'\underline{\underline{D}}^{-1}\underline{\underline{L}}\underline{\underline{L}}' + \underline{\underline{L}}' = \underline{\underline{L}}'\underline{\underline{D}}^{-1}\underline{\underline{S}}$$

or

$$(\underline{\underline{L}}'\underline{\underline{D}}^{-1}\underline{\underline{L}} + \underline{\underline{I}})\underline{\underline{L}}' = \underline{\underline{L}}'\underline{\underline{D}}^{-1}\underline{\underline{S}} \quad .$$

This and

$$\sum l_{rj}^2 + d_r = s_{rr} \text{ for } r=1, \dots, p$$

of (A6) are the maximum likelihood equations that have to be solved for  $l_{rj}$  and  $d_r$  .

References

- Armstrong, J. Scott (1967). Derivation of theory by means of factor analysis, or Tom Swift and his electric factor analysis machine. The American Statistician 21, 17-21.
- Harmon, H. H. (1967). Modern Factor Analysis, 2nd Ed. Chicago University Press.
- Joreskog, K. G. (1967). Some contributions to maximum likelihood factor analysis. Psychometrika 32, 443-482.
- Joreskog, K. G. and van Thillo, M. (1971). New rapid algorithms for factor analysis by unweighted least squares, generalised least squares and maximum likelihood. Research Memorandum, Educational Testing Service, Princeton, N. J. (May 1971)
- Kendall, M. G. and Lawley, D. N. (1956). The principles of factor analysis. J. Roy. Stat. Soc. A, 119, 83-84.
- Lawley, D. N. and Maxwell, A. E. (1971). Factor Analysis as a Statistical Method, 2nd Ed. American Elsevier Publishing Co., N. Y. (Review, JASA, March 1973, p. 224.)
- Press, S. J. (1972). Applied Multivariate Analysis. Holt, Rinehart and Winston, Inc., N. Y.
- Searle, S. R. (1971). Linear Models. Wiley, N. Y.