Department of Statistics P.O.Box 513 S-751 20 Uppsala Sweden

# SOME RESULTS ON USING SUMMED RAW SCORES AND FACTOR SCORES FROM DICHOTOMOUS ITEMS IN THE ESTIMATION OF STRUCTURAL EQUATION MODELS

Bengt Muthén April 1977

I am obliged to Bengt Dahlqvist for skillful computational assistance. Both his work and the author's was supported by the Bank of Sweden Tercentenary Foundation under project "Structural Equation Models in the Social Sciences", project director Karl G Jöreskog.

#### 1. Introduction

Statistical analyses in the social sciences frequently employ dichotomous variables that are measured in order to capture, in a more or less explicit way, some underlying or latent, quantitative variable. A typical example is a social-psychological study of attitudes or personality traits. Here, the response is often in the form of categorized answers to questionnaire statements. Another example is found in educational testing, with responses to items measuring a certain ability.

A common practice in such applications is to sum the scores (often set to 0 and 1) of the dichotomous items, forming a "scale", or an "index". This composite is assumed to have sufficient scale properties, so that it can be used in subsequent analyses assuming interval or ratio scales. In educational testing the number of items in such a scale is often larger than 50, while sociological analyses most often utilize small sets of items, even as few as 3 or 4. The present paper is mainly aiming at the latter types of applications.

In the major part of this paper, we are concerned with the following broad issue. How well do results from using a scale of this type agree with results from using the latent variable itself? Specifically, we will consider this question in the context of a structural equation model in which we have a system of linear relations between the latent variables and other variables involved. The measurement relations between the latent variables and the latent variables and the dichotomous items are assumed to fulfill the assumptions of the factor analysis model for dicho-

tomous variables (see Bock & Lieberman, 1970; Christoffersson, 1975; Muthén, 1977a). This is related to traditional latent trait models of test theory (see e.g. Lord & Novick, 1968;) see also Muthén (1977a). The structural equation model is a generalization of a model put forward by Muthén (1976), who treated estimation by the maximum-likelihood method. This estimator has desirable statistical properties, but as it turns out it is computationally extremely heavy. Simpler estimation procedures must be sought. Using summed raw scores of the dichotomous items, presents one such very simple estimation procedure. We will also study another comparatively simple estimator, which uses somewhat more response information from items than merely their sum. This uses factor scores, the i.e. we obtain an estimate of a certain latent variable for each observation in the sample.

The aim of this paper is to study the consistency of the two simple estimation procedures (described in more detail below). This may be seen as relating them to the "optimal", but computationally unfeasible maximum-likelihood estimator. We will not consider the variance of the estimators. These are hard to obtain for any general case of the model. Furthermore, the paper is particularly concerned with applications involving a small number of items. In such cases, it turns out that the asymptotic bias of the estimators is so large that questions of variance are comparatively uninteresting.

### A general model

.

Throughout the paper we will use the notation  $\mathbb{Z}_{wz}$  for the covariance matrix of the vectors of variables w and z. Also,  $\varepsilon_{wz}$  and  $\varepsilon_{wz}$  are to be understood in an analogous.

way. The variance of w is denoted  $\sigma_{ww}$ . Denote by u a p-dimensional vector of dichotomous items and by z a q-dimensional vector of quantitative (interval or ratio scaled) variables. These observed variables will now be related to the latent variables  $\eta(m \times 1)$  and  $\xi(n \times 1)$ , the relations between which are of primary interest. Each latent variable is measured either by variables from the u-vector or by variables from the z-vector. Denote by  $v(k \times 1)$  the vector of  $\eta$ - and  $\xi$ -variables of the first type and by  $u(q \times 1)$  the vector of  $\eta$ - and  $\xi$ -variables of the second type (here m + n = k + q).

Using an intervening p-dimensional vector v, v, v may be related to v as

$$u_{i} = \begin{bmatrix} 1, & if & v_{i} & \geq \tau_{i} \\ u_{i} & \geq & \tau_{i} \\ 0, & if & v_{i} & \leq \tau_{i} \end{bmatrix}$$
 (2.1)

i = 1, 2, ..., p, with

$$v = \Lambda v + \varepsilon . \tag{2.2}$$

Here,  $\tau$ (p x l) and  $\Lambda$ (p x k) contain measurement parameters. We assume that  $\nu$  and the residual vector  $\varepsilon$  are uncorrelated and multivariate normally distributed with zero means, and that  $\varepsilon$  has the diagonal covariance matrix

$$0 = I - \operatorname{diag}(\Lambda \Sigma_{\nu\nu} \Lambda^{-}) , \qquad (2.3)$$

where  $\sum_{\nu\nu}$  is the covariance matrix of  $\nu$ . This measurement specification is also used in the factor analysis model for dichotomous variables; for a fuller description the reader is referred to Bock (1970), Christoffersson (1975), and Muthén (1977a).

For simplicity, we will assume that  $\mu$  is measured without error by z

$$z \equiv \mu$$
 (2.4)

and that z has zero expectation (the results presented below are applicable also when allowing multiple indicators and measurement errors as in Jöreskog 1973, 1976).

For the latent variables under study, we will assume the linear structural equation system

$$B\eta = \Gamma\xi + \zeta , \qquad (2.5)$$

where B is a m x m non-singular parameter matrix,  $\Gamma$  is a m x n parameter matrix and  $\zeta$  is a m-dimensional disturbance vector, that is uncorrelated with  $\xi$  and has zero expectation. The covariance matrices of  $\xi$  and  $\zeta$  are denoted  $\Phi$  and  $\Psi$  respectively. We will assume that the vector  $(n^-, \xi^-)$  is multivariate normal with zero expectation. Thus,  $(v^-, z^-)$  is multivariate normal with mean vector zero and covariance matrix

$$\begin{bmatrix} \Sigma_{\nu\nu} & \Sigma_{\nuz} \\ \Sigma_{z\nu} & \Sigma_{zz} \end{bmatrix}, \qquad (2.6)$$

say.

The arrays  $\tau$ ,  $\Lambda$ , B,  $\Gamma$ ,  $\Phi$ , and  $\Psi$  contain the parameters of the general model. In a given application some of these parameters will have to be constrained in order to make the model identified.

The special case where v = n and  $u = \xi$ , i.e. where the dichotomous variables only appear as indicators of endo-

genous variables in (2.5), was treated by Muthén (1976).

## The two estimators

We will consider the estimation of the parameters of the general model, included in the structural relations of (2.5): B,  $\Gamma$ ,  $\Phi$ , and  $\Psi$ . In this section we define two comparatively simple estimation procedures for some special cases of the general model. The estimators will only be defined for models where  $\tau$ ,  $\Lambda$  and  $\Sigma_{\nu\nu}$  are identified from the marginal distribution of  $\Psi$ . The pxp covariance matrix of  $\nu$  is (see (2.2))

$$\sum_{n=1}^{\infty} = \sum_{n=1}^{\infty} \sum_{n=1}^{\infty} \sum_{n=1}^{\infty} \frac{1}{n}$$
 (3.1)

where diag  $(\Sigma)$  = I. To determine the metric of  $\nu$  we can either fix one element in each column of  $\Lambda$  to a non-zero value or fix the diagonal elements of  $\Sigma_{\nu\nu}$ .

For both estimators some knowledge is required about  $\tau$ ,  $\Lambda$  and  $\Sigma_{\nu\nu}$ . For the summed raw score estimator we must know the pattern of loadings in  $\Lambda$ , although not their exact values. It is to be assumed that the items of u have been carefully selected to measure  $\nu$ , so that  $\Lambda$  has a simple structure with several zero elements. For the factor score estimator it is necessary to know all values of  $\tau$ ,  $\Lambda$  and  $\Sigma_{\nu\nu}$ . If not known from previous analysis, these parameters can be estimated from the sample. This estimation may be carried out by the generalized least-squares factor analysis method of Muthén (1977a). In this case it will be assumed that the sampling errors can be ignored, and the parameters will be treated as given in subsequent analysis.

In this context we should note the special case of the general model, where there are no z's involved (i.e. k = m+n). Then we can directly use  $\Sigma_{\nu\nu}$ , as estimated from the factor analysis, in the second step of the estimation. (If in fact  $\Sigma_{\nu\nu}$  is unrestricted, the structural estimates are directly obtained, without loss of efficiency in the second step).

For both estimators, the first step of the estimation procedure results in a sample of vectors that are proxies to the sample of  $\nu$ -vectors (although in a different metric). In their common second step, the second-order sample moments of these proxies and of z are used to estimate the parameters of (2.5). For simplicity we will here assume that these sample moments give a sample covariance matrix ([k+q] x [k+q]) that is positive definite. The computation of the estimates of  $B, \Gamma, \phi$ , and  $\Psi$  from this covariance matrix is then straight-forward, and may for instance be carried out by the maximum-likelihood method of the general LISREL procedure (see Jöreskog, 1973, 1976). When the parameters of (2.5) are just-identified in terms of the covariance matrix of ( $\nu$ , z), i.e. when this matrix is unrestricted, the maximum-likelihood estimator can be given explicitly. This will be utilized in some examples below.

# 3.1. The summed raw score estimator

For this estimator it will be convenient to define a p x k weight matrix W, such that

$$\underline{y} = \underline{W} \underline{u}, \qquad (3.2)$$

where  $\underline{y}$  is the k-dimensional vector of summed raw scores. When all rows of  $\underline{\Lambda}$  have only one non-zero element, the items of  $\underline{u}$  will be called "pure indicators" of their respective latent variable (corresponding to the columns of  $\underline{\Lambda}$ ). We may then assume that the  $\underline{u}$ -alternatives are labelled so that all elements

of  $\Lambda$  are non-negative, and define the i,j-th element of W as

$$[W]_{ij} = \begin{cases} 1, & \text{if } \lambda_{ij} \neq 0 \\ 0, & \text{if } \lambda_{ij} = 0 \end{cases}$$
 (3.3)

When some indicators (u-variables) are directly related to more than one latent variable the choice of W is not always evident. In this paper we will only consider examples where all elements of  $\Lambda$  are non-negative (at least after interchanging some u-alternatives), i.e. all elements in a certain row of  $\Lambda$  have the same sign. In this case, it is still relevant to define W as in (3.3). However, the derivations that follow are quite general in that any  $\Lambda$  can be used together with any constant matrix W.

From a sample of u-vectors, the transformation (3.2) gives a sample of y-vectors which may be used in the second step of the estimation procedure discussed above.

## 3.2. The factor score estimator

We will consider the estimation of  $\nu$  in (2.2) by a factor score estimator of a Bayesian type, suggested by Samejima (1969) for the case of k=1, and generalized to the multiple factor case by Muthén (1977b). This estimator requires that the values of  $\tau$ ,  $\Lambda$ , and  $\Sigma_{\nu\nu}$  are given. The factor score estimator maximizes the density of the distribution of  $\nu$  conditional on  $\nu$ , with respect to the elements of  $\nu$ . We will denote the estimated k-dimensional factor score vector by f. From a sample of  $\nu$  swe will thus obtain a sample of  $\nu$  so In the second step of our estimation procedure, we note that there is a choice of using  $\Sigma_{\nu\nu}$  as estimated in the factor analysis (assuming that  $\Sigma_{\nu\nu}$  is not known), or using the sample covariance matrix of f. In this paper we will only study the factor score estimator for the

case of k = 1, determining the metric of  $\nu$  through the stand- ardization  $\sigma_{\nu\nu}$  = 1, and we will therefore not adress this issue.

An alternative factor score estimator is the traditional conditional maximum-likelihood method, see e.g. Lord (1968), Samejima (1974), Muthén (1977b). However, this estimator (assuming k = 1 and loadings of equal sign) is not suitable since it yields infinite factor score estimates for the "extreme" u-response patterns (1 1...1) and (0 0...0). We are particularly concerned with situations involving a small number of items. Then these extreme patterns are likely to have a non-negligible number of observations, which cannot be used in the second step of the estimation procedure. We finally note that in the case of quantitative response variables, the Bayesian factor score estimation approach results in the ordinary regression method with correlated factors (see e.g. Lawley & Maxwell, 1971).

# 4.1. The bias of the summed raw score estimator

We will first determine  $\sum_{yy} (k \ x \ k)$  and  $\sum_{yz} (k \ x \ q)$ . The i,j-th element of  $\sum_{uu} (p \ x \ p)$  is

$$P(u_{i} = 1) \cdot [1 - P(u_{i} = 1)], \text{ if } i = j$$

$$\{ \sum_{uu} \}_{ij} = P(u_{i} = 1, u_{j} = 1) - P(u_{i} = 1) \cdot P(u_{j} = 1)$$

$$\text{if } i \neq j.$$

$$(4.1)$$

From the model of Section 2 we find that

$$P(u_i = 1) = \int_{\tau_i}^{\infty} \phi(z) dz$$
 (4.2)

and for i # j

$$P(u_{i} = 1, u_{j} = 1) = \int_{\tau_{i}\tau_{i}}^{\infty \infty} \phi(z; 0, \Sigma_{ij})dz$$
. (4.3)

Here  $\phi(z)$  denotes the density of the univariate, standardized normal distribution. We use  $\phi(z; a, B)$  as the notation for the density of a multivariate normal distribution with mean vector a and covariance matrix b. In (4.16) the density corresponds to a bivariate distribution where the covariance matrix  $\sum_{ij}$  is formed from (3.1) with  $\left[\sum_{ij}\right]_{ii}$  and  $\left[\sum_{ij}\right]_{jj}$  as diagonal elements and  $\left[\sum_{ij}\right]_{ij}$  as off-diagonal elements. In Section 4.1 we will consider cases where  $\lambda=1$  for k=1 and  $\sigma_{\nu\nu}=1$ . Then  $\lambda_i=\lambda_j=1$  gives  $\left[\sum_{ij}=1$ . In such rare cases  $\nu_i=\nu_j^*=\nu$ , and  $P(u_i=1,u_j=1)=P(\nu\geq\max_i \{\tau_i,\tau_j\})=\min_{ij}\left[P(u_i=1),P(u_j=1)\right]$ . We have now determined  $\sum_{uu}$ , giving

$$\sum_{\mathbf{y}\mathbf{y}} = \mathbf{W}^{\mathbf{x}} \sum_{\mathbf{u}\mathbf{u}} \mathbf{W} . \tag{4.4}$$

To determine  $\Sigma_{yz}$  we note that <u>conditional</u> on  $\nu$  the model gives

$$\sum_{\mathbf{y}\mathbf{z}} = 0 \qquad (4.5)$$

Then

$$E(y \cdot z^-|v) = E(y|v) \cdot E(z^-|v), \qquad (4.6)$$

where

$$E(z^-|v) = v^- \sum_{v=v}^{-1} \sum_{v=v} (4.7)$$

Put

$$E(u|v) = \pi$$
 (4.8)

The i-th element of  $\pi$ ,

$$\pi_{i} = P(u_{i} = 1 | v) = \int_{\tau_{i}}^{\infty} \phi(z; \lambda_{i}, v, \theta_{ii}^{-1/2}) dz$$
, (4.9)

where  $\lambda_{i}$  is the i-th row of  $\Lambda$  and  $\theta_{i,i}$  is the i-th diagonal element of  $\Theta$ . Now,

$$\sum_{yz} = E(\underline{y} \cdot \underline{z}) =$$

$$= E[E(\underline{y} \cdot \underline{z} | \underline{v})], \qquad (4.10)$$

where E denotes the expectation with respect to the (marginal) v distribution of v. As is seen from (4.10) and (4.6) it remains to determine

$$E[E(y|v) \cdot v] = W E(\pi \cdot v). \tag{4.11}$$

In Appendix it is shown that

$$E(\pi \cdot \nu) = D_{\phi} \Lambda \Sigma_{\nu\nu} , \qquad (4.12)$$

where D is a diagonal matrix with  $\, \varphi \, (\tau_{\, 1}^{}) \,$  as i-th diagonal element. It follows that

$$\sum_{\mathbf{y}z} = \mathbf{W}^{\mathbf{D}} \mathbf{D}_{\phi} \mathbf{\Lambda} \mathbf{\Sigma}_{\mathbf{v}z} . \tag{4.13}$$

To be able to compare the population moments of the summed raw scores with those of  $\nu$ , we will transform  $\nu$  to the same metric as  $\nu$ , creating

$$y^* = p_v^{1/2} p_y^{-1/2} y, \qquad (4.14)$$

with diag( $\sum_{y=y}^{x}$ \*) = diag( $\sum_{y y}$ ). Putting

$$D_{\nu}^{1/2}D_{y}^{-1/2}WD_{\phi}\Lambda = C ,$$
it follows by (4.13) that
$$C = [WD_{\phi}\Delta]^{-1/2}D_{\phi}\Delta = C ,$$

$$C = [WD_{\phi}\Delta]^{-1/2}D_{\phi}\Delta = D_{\phi}\Delta = D_{\phi}\Delta$$

$$\sum_{z} *_{z} = C\sum_{z > 0} . \tag{4.16}$$

In passing we note that the k xk matrix of correlations

between y and v is

$$\Omega_{yv} = D_{y}^{-1/2} W^{-} D_{\phi} \Lambda \Sigma_{vv} D_{v}^{-1/2} 
= D_{v}^{-1/2} C \Sigma_{vv} D_{v}^{-1/2} ,$$
(4.17)

Generally,  $\sum_{y} *_{y} *_{z} \neq \sum_{vv}$  and  $\sum_{y} *_{z} \neq \sum_{vz}$ , i.e.  $C \neq I$ . This means that any consistent estimator of  $\sum_{yv}$  and  $\sum_{yz}$  will not be a consistent estimator of  $\sum_{vv}$  and  $\sum_{vz}$ , taking the difference in metric into account. Thus, using the sample moments of y and z will generally give asymptotically biased estimates of the parameters of (2.5): B,  $\Gamma$ ,  $\Phi$ , and  $\Psi$ . Through some examples presented below, we will try to give a picture of the size of this bias.

$$c = \rho , \qquad (4.18)$$

with  $\rho$  denoting the correlation between y and  $\nu$ ,

$$\rho = \sigma_{vv}^{1/2} \sigma_{yy}^{-1/2} \sum_{i \in S} \phi(\tau_i) \lambda_i, \qquad (4.19)$$

where S is the set of u's considered. If all rows of  $\Lambda$  have only one non-zero element, C is a diagonal matrix with the different  $\rho$ 's on its diagonal. The formula for  $\rho$  in the special case of equivalent items, was given by Tucker (1946) in a study of test validity. We conclude that a consistent estimator of the covariances of y with the other variables in the structural relations of (2.5) will in this case be biased down-ward with a factor  $1-\rho$  (0 <  $\rho$  < 1) relative to the corresponding covariances of  $\nu$ .

#### 4.1 Some examples

A computer routine has been created for the calculation of  $\sum_{x,y}^{*}$  and  $\sum_{x}^{*}$  and  $\sum_{x}^{*}$  and  $\sum_{x}^{*}$  for the computation of (4.3), this routine uses parts of an algorithm developed by Kirk (1973).

Let us consider the case of k=1, or similarly the case of a certain latent variable having pure indicators. We will first show how  $\rho$  varies with the number of dichotomous items and with the  $\lambda$ -parameters. For simplicity we will denote the number of items by  $\rho$ . In Figure 1,  $\rho$  is given for equivalent items with  $\tau=0$  and  $\lambda$  varying from .2 to 1, and for a common range of p-values (with a different parameterization, a similar figure was given in Tucker, 1946, including  $\rho=100$ ). We have used  $\sigma_{\nu\nu}=1$ .

INSERT FIGURE 1 ABOUT HERE

\_\_\_\_\_\_

Social - psychological applications frequently have p-values around 5 and  $\lambda$ -values around .6, which by the figure would result in a downward covariance bias of about 23%.

Consider for instance the case where k = 1 and the v-variable appears as the only endogenous variable of (2.5)

$$v = \gamma_1 z_1 + \gamma_2 z_2 + \dots + \gamma_q z_q + \zeta. \tag{4.20}$$

Here,

$$Y = \sum_{z=z}^{-1} \sigma_{zy} , \qquad (4.21)$$

where  $\gamma' = (\gamma_1, \gamma_2, ..., \gamma_q)$ . In the regression of  $y^*$  on z, we have (see (4.16))

$$\chi^* = \rho \cdot \chi . \tag{4.22}$$

With 5 indicators of  $\nu$  and  $\tau = 0$ ,  $\lambda' = (6, .6, ..., .6)$ ,  $\sigma_{\nu\nu} = 1$ , the downward asymptotic bias in a consistent estimator of the regression coefficients will thus be 23%. The portion of "explained" variance in  $\nu$ ,  $\gamma' \Sigma_{zz} \gamma$ , will be estimated to an approximation of only 59% (i.e.  $\rho^2$ ) of its true value.

From Figure 1 we see that for given  $\tau$ -values,  $\rho$  is an increasing function of the  $\lambda$ 's throughout the common range of  $\lambda$ . The fact that  $\rho$  decreases for  $\lambda$ 's near 1 has been called "the attenuation paradox"; see e.g. Lord & Novick (1968, p 344). If we accept 10 - 15% or less as a negligible size of covariance bias for the model of (4.20), we note that we still need around 15 items or more for the most common range of  $\lambda$ , say  $0.4 \le \lambda \le 0.8$ .

indication

We will now give an identification of how  $\rho$  varies with  $\tau$ . We recall that large positive (negative) values of  $\tau$  are associated with items for which the probability of observing u = 1 is small (large). We will use equal loadings set to .6, which will be employed as a standard value. First,

consider p = 5. With  $\tau = 0$ , Figure 1 gives  $\rho = .769$  (as was found above). As a standard example of unequal  $\tau$  s, these will be evenly placed along the v-scale, in the sense that the area between the corresponding p + 1 intersections of the normal curve are of equal size. Thus, with  $\tau = (-1.00, -.50,$ .00, .50, 1.00) we obtain  $\rho = .741$ . To illustrate cases of  $\tau$ 's covering only one end of the  $\nu$ -scale, we may use  $\tau$ ' = (.00, .50, .50, 1.00, 1.00). Here p = .730. The examples indicate that  $\rho$  decreases with increasing distance of the  $\tau$  s from the zero mean of v. Consider the case of a somewhat larger number of items, say p = 8. Here,  $\tau = 0$  gives  $\rho = .833$ . With evenly distributed  $\tau$ 's in the sense used above, we have  $\tau = (-1.25, -.75, -.45, -.15, .15, .45, .75, 1.25)$  and  $\rho = .809$ . With  $\tau = (.15, .15, .45, .45, .75, .75, 1.25, 1.25) \rho = .817. We$ note that in all cases the p-values differ only slightly from the value obtained with  $\tau = 0$ .

As an example with different  $\lambda$ 's, consider the case of p=5,  $\tau=0$ , and  $\lambda'=(.4,.5,.6,.7,.8)$ . Here,  $\rho=.770$ , i.e. about the same as with equal  $\lambda$ 's with the value .6. For p=8,  $\tau=0$  and  $\lambda'=(.4,.5,.5,.6,.6,.7,.7,.8)$  we have  $\rho=.833$ , as in the previous example with equal  $\lambda$ 's. As a slight digression, let us consider using different weights in W, with weights set equal to the loadings. For the previous two examples we then obtain  $\rho=.788$  and  $\rho=.843$ , respectively. Furthermore, if we have 5 items with zero  $\tau$ 's and equal loadings of .6, and add 5 items with zero  $\tau$ 's and equal loadings of .3,  $\rho$  only increases from .769 to .773. Weighting with the loadings gives an increase to .797.

This collection of examples also serve the purpose of de-

monstrating the usefulness of Figure 1. We see that the  $\rho$ -values obtained from this figure can be used as good approximations in cases where we do not have equivalent items with  $\tau = 0$ , and where we use somewhat different weights.

Let us now turn to some examples where k = 2. Assume that 12 items are available to measure the two latent variables and that the first 10 of these load with .6 on only one of the latent variables, 5 items on each variable. The remaining 2 items load equally on the latent variables, and we will use the loading .35 as the common value, so that the residual variances of  $\Theta$  are about equal for all items. Let us further assume that  $\tau = 0$  and that we have standardized the variances of the latent variables to one, with a correlation of .5. Given this, we can study the difference between  $\Sigma_{yy}$  and  $\Sigma_{yy}$ , and between C and I. Using the first 10 items we obtain

$$\sum_{\mathbf{Y}} \mathbf{Y} \mathbf{Y} = \begin{bmatrix} 1.000 \\ \\ .297 \\ 1.000 \end{bmatrix} \qquad \sum_{\mathbf{C}} = \begin{bmatrix} .769 \\ \\ .000 \\ .769 \end{bmatrix} (4.23)$$

and using all 12 items,

$$\Sigma_{\mathbf{Y}^*\mathbf{Y}^*} = \begin{bmatrix} 1.000 \\ .581 & 1.000 \end{bmatrix} \quad \tilde{\mathbb{C}} = \begin{bmatrix} .732 & .139 \\ & & \\ .139 & .732 \end{bmatrix}.$$

We can now pose the question: To obtain the smallest bias in the structural coefficients, should we utilize only the first 10 items, obtaining two sets of pure indicators, or should we make use of all 12 items? In fact, the answer depends

on the structural model. Consider the following two simple models. In Model 1,

$$z = [\gamma, \gamma] \nu + \zeta$$
, (4.25)

and in Model 2

$$v = \begin{bmatrix} v \\ y \end{bmatrix} z + \xi . \tag{4.26}$$

In both cases we standardize the variance of z to one. Denote by  $\gamma^*$  the vector corresponding to  $[\gamma, \gamma]$ , when using  $\gamma^*$  instead of  $\gamma$ . For the first model we have, see (4.16),

$$y^* = \sum_{y=y}^{-1} {}^*C\sigma_{yz}$$
, (4.27)

and for the second one,

$$\gamma = C\sigma_{VZ} \quad . \tag{4.28}$$

This implies that for Model 1 the summed raw score estimator will give a down-ward bias of 11% in the estimation of  $\gamma$ , when using only the 10 items. For Model 2 these items give a downward bias of 23% (as was obtained for the model of (4.20)). Using all 12 items, there is a 17% down-ward bias for Model 1 and a 13% down-ward bias for Model 2. Thus, the relation between the sizes of bias is reversed between these two models.

#### 4.2. Correcting for bias

In the previous examples we have found the bias in the estimates of the structural coefficients, using y instead of v, as different functions of the elements of v, and v, are setting elements of v, and v, and v, and v, and v, are setting elements of v, and v, and v, and v, and v, are setting elements of v, and v, and v, and v, are setting elements of v, and v, are setting elements of v, and v, are setting elements of v, and v, and v, are setting elements of v, are setting elements of v, and v, are setting elements of v, are setting elements of v, and v, are setting elements of v, are setting ele

Muthén (1977a). This suggests that we may use these estimates to correct for the bias in the sample covariance matrices  $S_{yy}$  and  $S_{yz}$ .

Consider the corrected sample covariance matrix

$$\mathbf{S}_{yz}^{*} = \hat{\mathbf{C}}^{-1}\mathbf{S}_{yz}^{-1}, \quad \hat{\mathbf{D}}_{y}^{1/2}\hat{\mathbf{O}}_{y}^{-1/2}\hat{\mathbf{C}}^{-1}\hat{\mathbf{O}}_{y}^{1/2}\hat{\mathbf{O}}_{y}^{-1/2}\hat{\mathbf{O}}_$$

where we have assumed that the correction matrix  $\hat{\mathbb{C}}$  is non-singular (this will always be the case for pure indicators, since then  $\hat{\mathbb{C}}$  is diagonal with the different  $\hat{\rho}$ 's on its diagonal). Here,  $\hat{\mathbb{C}}$  is created from the estimated  $\tau$ ,  $\Lambda$ , and  $\Sigma_{\nu\nu}$ . From (4.15) and (4.13) it follows that  $\Sigma_{\nu\nu}^*$  is a consistent estimator of  $\Sigma_{\nu\nu}$ , when transformed to the metric of  $\nu$ . (In actual practice, we would usually not change the metric of  $\nu$ , i.e. we would directly use  $\Sigma_{\nu\nu}^*$ ).

When k > 1, we must also take the bias in the covariances of y into account. We can then use  $\hat{\Sigma}_{\nu\nu}$ . With  $\hat{S}_{\nu\nu}^*$  and  $\hat{\Sigma}_{\nu\nu}$  in the same metric, and adding  $\hat{S}_{zz}$ , we hopefully obtain a positive definite covariance matrix, which is a consistent estimator of the covariance matrix of  $(\hat{\nu}, \hat{z})$ . This may then be used to obtain a consistent estimator of the structural parameters.

## 5. The bias of the factor score estimator

In this section we will compare the bias of the summed raw score estimator with that of the factor score estimator. The estimator of the factor scores is defined as the minimum of a certain function (see Samejima, 1969; Muthén, 1977b), and cannot be given explicitly. Due to this, we will investigate its properties by means of a small Monte Carlo

study. This will be limited to the case of one latent variable. Analogous to Section 4, we will mainly concentrate on the covariance between the latent variable  $\nu$  and another variable in the model, z. Thus, we are particularly interested in the bias of  $\sigma_{\rm fz}$ . We will also report results pertaining to  $\mu_{\rm f}$ , i.e. the mean of f,  $\sigma_{\rm ff}$ , and the regression of z on f,

$$\gamma_{zf} = \sigma_{ff}^{-1} \sigma_{fz} . \qquad (5.1)$$

The Monte Carlo study was designed in the following way. Given certain values of  $\sigma_{zz}$ ,  $\sigma_{vz}$ ,  $\sigma_{vv}$ , p,  $\tau$ , and  $\lambda$  (a column of  $\Lambda$ ) the (p+1) x (p+1) population covariance matrix of  $(v^*,z)$  was created; see Sections 2 and 3. We used the values  $\sigma_{zz}=1$ ,  $\sigma_{vz}=.5$ , and  $\sigma_{vv}=1$ ; thus, in this case  $\sigma_{vz}$  is a correlation coefficient. From the general model we find that the  $v^*$ s have unit variances, the covariance of  $v^*_i$  and  $v^*_j$  is  $\lambda_i \lambda_j$ , and the covariance between  $v^*_i$  and z is  $\lambda_i \sigma_{vz}$ . Using this covariance matrix, 50 samples of 2000 random, multivariate normal vectors were created. The first p variables of each random vector, corresponding to the  $v^*$ s, were dichotomized at  $\tau$ , creating u as in (2.1). Given the pattern of u, and given  $\tau$  and  $\lambda$ , t was computed with the Bayes factor score estimator. This resulted in a pair of t-and t-values for each random observation unit.

We will use a "curl"-notation for the means of the different characteristics, taken over the 50 values. Thus,  $\sigma_{\rm ff}$  denotes the mean of the 50 values of  $s_{\rm ff}$ . As it turns out,  $\sigma_{\rm ff}$  differs from  $\sigma_{\rm VV}$ . Since the metric of the latent variable is arbitrary, this is no deficiency. However, as in Section 4,

we will use a transformation

$$f^* = \tilde{\sigma}_{ff}^{-1/2} f, \qquad (5.2)$$

so that  $f^*$  and v are (approximately) in the same metric. Then  $\sigma_{fz}^* = \tilde{\sigma}_{ff}^{-1/2} \tilde{\sigma}_{fz}$ , an approximation to the correlation between f and z. The standard error of  $\sigma_{fz}^*$  is equal to  $\tilde{\sigma}_{ff}^{-1/2}$  multiplied by the standard error of  $\sigma_{fz}$  (i.e.  $\sigma_{ff}$  is here treated as a constant, approximating  $\sigma_{ff}$ ).

Let us now consider some of the examples that were used in connection with the summed raw score estimator. We will limit the study to three examples with p = 5 and three examples with p = 8. It is for p-values of this magnitude that it is possible to obtain any substantial reduction of bias relative to using the summed raw scores. As a check of the precision of the Monte Carlo procedure, the summed raw scores were also calculated from the dichotomization of v. In all the examples reported, the calculated  $\sigma_{v}$  differed less than .004 from  $\sigma_{v}$ . In Table 1 the results are displayed. We recall that  $\tilde{\sigma}_{f}$  shall be compared to  $\sigma_{vz}$  = .5.

INSERT TABLE 1 ABOUT HERE

It is somewhat surprising to note that in none of the examples does the factor score estimator perform distinctly better than the summed raw score estimator. The additional information used by this estimator, giving a considerable extra computational work, seems to give a very small difference in bias, if any at all. Judging from the sample standard errors of the  $\tilde{c}_f^*$ , only one of these values are significantly different from  $\rho \cdot c_{\sqrt{2}}$ 

(the sampling distributions of the sf\*z's are approximately normal). The largest differences occur for cases with unequal loadings.

We can note that  $\hat{\gamma}_{zf}$  is comparatively close to the correct value of .5, in all cases studied (note that  $\hat{\gamma}_{fz} = \hat{\sigma}_{f} *_{z}$ ). It appears that the bias in  $\sigma_{ff}$  to some extent balances the bias in  $\sigma_{fz}$ . It is interesting to note that this holds true exactly for the case of quantitative variables, using the related regression method (see Tucker, 1971). Also for this method, the covariances of the estimated factor scores with other variables differ from the corresponding covariances of the true factor scores (Tucker, 1971).

## Summary and discussion

We have shown that the use of summed raw scores to estimate structural equation parameters can give grossly biased results when the number of items is small. On the other hand, when each latent variable is measured by a larger number of items, say around fifteen or more, this procedure may work quite well with regard to bias. Due to its simplicity of calculation, the method of using summed scores will continue to be attractive. We have therefore attempted to give a picture of how the bias varies under different conditions, indicating when the bias is small.

With one latent variable, the bias of the factor score estimation procedure is close to that of using summed scores. Then the latter will probably be preferred on grounds of simplicity.

There is also the possibility of correcting the sample moments of the summed scores, as discussed in Section 4.2. In this correction we utilize parameter values that are also necessary for the calculation of the factor scores. Even with correction, the summed scores are simpler to calculate, given a routine of the type used in Section 4.1.

In this paper we have only considered the bias of the estimators. Thus, questions of standard errors for the estimates and statistical tests of model fit have not been discussed. Such information is important in the evaluation of a model. To conclude, we note that for the type of models discussed, there is still a need to develop estimators that are computationally feasible, and yet statistically acceptable. This is especially true for models with a small number of items. It is desirable that such an estimator also produces standard errors and a test of model fit.

Good way to go - Simple model

#### APPENDIX

In this appendix we will use the same notation as in the main part of the paper. Using standard results on conditional normal distributions, we find that

$$(A-1) \qquad \phi(z; \widetilde{Aw}, B) \cdot \phi(w; e, F) =$$

$$= \phi(z; Ae, B + AFA^{-}) \cdot \phi(w; \mu_{w \cdot z}, \Sigma_{w \cdot z}),$$

where

(A-2) 
$$\mu_{w \cdot z} = e + FA^{-}(B + AFA^{-})^{-1} \cdot (z - Ae) ,$$

$$\Sigma_{\mathbf{W}^{\bullet}\mathbf{Z}} = \mathbf{F} - \mathbf{F}\mathbf{A}^{-}(\mathbf{B} + \mathbf{A}\mathbf{F}\mathbf{A}^{-})^{-1} \mathbf{A}\mathbf{F} .$$

Denote by

$$\prod_{\infty}^{\infty} \tilde{t}(\tilde{m}) q\tilde{m}$$

a multiple integral which has the order of the dimension of w, upper and lower limits all equal to positive and negative infinity, respectively, and the vector-valued function f(w) as integrand. Consider the i-th row (i = 1, 2, ..., p) of  $E(\pi \cdot v)$  of (4.11):

$$E(\pi_{i} \cdot v) = \int_{-\infty}^{\infty} \int_{\tau_{i}} \phi(z; \lambda_{i}, v, \theta_{ii}) dz \cdot v.$$

$$(A-4) \qquad \qquad \cdot \phi(\nu; 0, \Sigma_{\nu\nu}) d\nu.$$

Changing the order of integration between z and v, and using (A-1), (A-2), (A-3), we find

(A-5) 
$$E(\pi_{1} \cdot \nu) = \int_{\tau_{1}}^{\infty} \varphi(z) \cdot z \cdot \lambda_{1}^{r} \xi_{\nu\nu} dz ,$$

where we have made use of the fact that  $\theta_{ii} + \lambda_{i} \sum_{\nu \nu} \lambda_{i} = 1$ , for all i's. Utilizing the expression for the mean of a truncated normal distribution (see e.g. Tallis, 1961), we obtain

(A-6) 
$$\underbrace{\mathbf{E}(\pi_{\mathbf{i}} \cdot \mathbf{v})}_{\mathbf{v}} = \phi(\tau_{\mathbf{i}}) \cdot \lambda_{\mathbf{i}} \cdot \sum_{\mathbf{v} \mathbf{v}}.$$

This gives the result of (4.12).

#### REFERENCES

- Bock, R.D. & Lieberman, M. Fitting a response model for n dichotomously scored items. Psychometrika, 1970, 35, 179 197.
- Christoffersson, A. Factor analysis of dichotomized variables. <u>Psychometrika</u>, 1975, <u>40</u>, 5 32.
- Jöreskog, K.G. Structural equation models in the social sciences: Specification, estimation and testing. Research Report 76-9. Department of Statistics, University of Uppsala, 1976.
- Jöreskog, K.G. A general method for estimating a linear structural equation system. In A.S. Goldberger and O.D. Duncan (eds.): Structural equation models in the social sciences. New York: Seminar Press, 1973, 85 112.
- Kirk, D.B. On the numerical approximation of the bivariate normal (tetrachoric) correlation coefficient. Psychometrika, 1973, 38, 259 -268.
- Lawley, D.N. & Maxwell, A.E. Factor analysis as a statistical method.

  London: Butterworths, 1971.
- Lord, F.M. An analysis of the verbal scholastic aptitude test using Birnbaum's three-parameter logistic model. Educational and Psychological Measurement, 1968, 28, 989 1020.
- Lord, F.M. & Novick, H. Statistical theories of mental test scores.

  Reading, Mass.: Addison-Wesley Publishing Co., 1968.
- Muthén, B. Structural equation models with dichotomous dependent variables. Research Report 76-17. Department of Statistics, University of Uppsala, 1976.
- Muthén, B. Contributions to factor analysis of dichotomous variables.

  Department of Statistics, University of Uppsala 1977a.
- Muthén, B. On the estimation of factor scores for dichotomous variables.

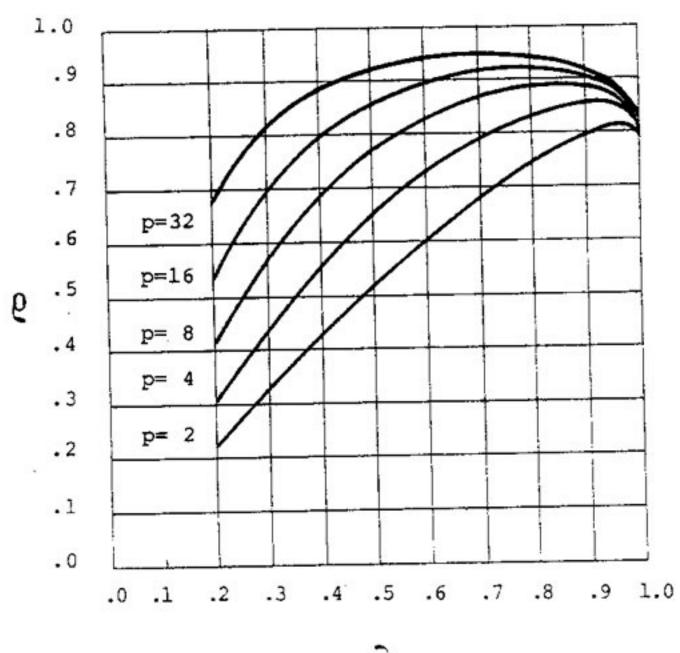
  Department of Statistics, University of Uppsala, 1977b.
- Samejima, F. Estimation of latent ability using a response pattern of graded scores. Psychometrika Monograph, 1969, No 17.
- Samejima, F. Normal ogive on the continous response level in the multidiemnsional latent space. Psychometrika, 1974, 39, 111 - 121.
- Tallis, G.M. The moment generating function of the truncated multinormal distribution. Journal of the Royal Statistical Society, Series B, 1961, 23, 223 229.
- Tucker, L.R. Maximum validity of a test with equivalent items. Psychometrika, 1946, 11, 1 - 13.
- Tucker, L.R. Relations of factor score estimates to their use. Psychometrika, 1971, 36, 427 - 436.

TABLE 1
Monte Carlo results for six sets of items.\*

μf	off	Yzf	of*z	ρ•υνΖ
= 5, τ	$= 0, \lambda^{-} = (.6,$			
.001	.535 (.002)	.531 (.004)	.388 (.003)	.385
σ = 5, τ	= (-1.0,5,	.0, .5, 1.0), %	= (.6, .6,	, .6):
.000	.511 (.002)	.520 (.004)	.372 (.003)	.370
= 5, τ	-	.5, .6, .7, .8	):	
	.525 (.001)	.550 (.004)	.399 (.003)	.385
= 8, τ	$= 0, \chi = (.6,$	.6,, .6):		
.003	.638 (.002)	.519 (.003)	.414 (.004)	.417
		5,45,15,	.15, .45, .75	, 1.25),
~	= (.6, .6,		10.00	
.003	.613 (.002)	.512 (.004)	.401 (.003)	.404
$\tau = 8$	$= 0, \lambda^- = (.4,$	.5, .5, .6, .	6, .7, .7, .8)	:
.001	.628	.531	.421	.417

<sup>\*</sup> Standard errors in parenthesis.

The correlation between the summed raw score and the latent variable.



λ