

Analysis of Qualitative Variables in Structural Models with Unique Solutions

Giorgio Vittadini

Università degli Studi di Milano,
20126 – Viale Sarca 202
Milano, Italia, e-mail: vittadin@imiucca.csi.unimi.it

Abstract: A new method based on the Multidimensional Scaling and the Restricted Regression Component Decomposition is proposed in order to obtain solutions for structural models with mixed variables.

Keywords: Structural Models with Mixed Variables, Normality Hypothesis, Multidimensional Scaling, Alternating Least Squares, Restricted Regression Component Decomposition.

1. Structural Model with Qualitative Variables

The structural model for the study of causal relationship among latent variables is composed of one structural and two measurement equations:

$$H = HB + \Xi\Gamma + E = \Xi\Gamma(I-B)^{-1} + E(I-B)^{-1}; Y = H\Lambda_y + U; X = \Xi\Lambda_x + U \quad (1)$$

where $Y' = (y'(1), \dots, y'(t))(n_y, t)$, $X' = (x'(1), \dots, x'(t))(n_x, t)$ are the observed mixed variables; $\Xi' = (\xi'(1), \dots, \xi'(t))(n_\xi, t)$, $H' = (\eta'(1), \dots, \eta'(t))(n_\eta, t)$ are the latent variables; $E' = (\varepsilon'(1), \dots, \varepsilon'(t))(n_\varepsilon, t)$ are the errors in equations; $\Delta' = (\delta'(1), \dots, \delta'(t))(n_\delta, t)$, $U' = (u'(1), \dots, u'(t))(n_u, t)$, are the errors in variables. It is assumed that: all the random variables have zero mean and finite variance, B is a low matrix with zero on the main diagonal, (Y, X, H) are identically distributed and (Ξ, E, Δ, U) are identically and independently distributed. The model is usually proposed with restrictions on parameters and on covariances.

The solutions are reached starting from the variance covariance matrix of the reduced model where the variables H in the measurement models are substituted by the value of $(\Xi\Gamma(I-B)^{-1} + (I-B)^{-1}E)$ obtained from the structural equations.

2. Methods for Obtaining Solutions

First of all given that a continuous bivariate normal variable (J_l^*, J_m^*) with distribution $\Phi(j_l^*, j_m^*, \rho_{(j_l^*, j_m^*)})$ underlies every pair of ordinal bivariate observed variables (J_l, J_m) , the polychoric correlation between the two components of them is calculated. When there is one only ordinal variable, the polyserial correlation coefficient between an observed quantitative variable and the normal underlying variable is defined. The observed frequencies of the qualitative variables n_{j_l, j_m} ($q = 1, \dots, n_l; r = 1, \dots, n_r$) given, the polychoric coefficients are reached in different ways. Jöreskog (1994) hypothesizes that the marginal probability of the normal variables are equal to the marginal frequencies of the two-way table ($\pi_{j_l, \cdot} = n_{j_l, \cdot}$) ($\pi_{\cdot, j_m} = n_{\cdot, j_m}$) of the ordinal variables. Therefore, first of all, he reaches the thresholds, using the marginal distributions of the normal bivariate distribution; for given thresholds using such distribution, he obtains the correlation coefficient maximizing the log likelihood of the sample respect to $\rho_{(j_l^*, j_m^*)}$. Lee et al. (1990) estimate the thresholds by means of Partition Maximum Likelihood (PML) which is simpler from the computational point of view. Lee et al. (1995) estimate simultaneously the correlation coefficients and the thresholds concerning pairs of variables maximizing the log likelihood of the sample respect to $\rho_{(j_l^*, j_m^*)}$ and respect to every threshold j_l^*, j_m^* , by means of PML (but Jöreskog (1994) observes that "different estimates of thresholds for one variable may be obtained from different pairs of variables"). By means of Full Maximum Likelihood, in one case Lee et al. (1992) simultaneously reach all the thresholds of the polychoric correlations; in another case (Lee et al. (1990)), they reach also the parameters, the variances and the covariances of the latent variables. Moreover these two last methods take up too much computer time (Lee et al. (1990) Lee et al. (1995)). In the first three methods, the parameters and the covariances of the latent variables and errors are obtained from the polychoric or polyserial correlations. There is not a unique understanding about the method of obtaining the parameters and the latent variables. Jöreskog (1990), (1994), Rigdon and Ferguson (1991) propose the Weighted Least Squares method and criticise the Maximum Likelihood method because the standard error parameter estimates are asymptotically incorrect, but Lee et al. (1995) continue to prefer the General Least Squares method and criticize the proposal of Jöreskog because it requires sample sizes larger than 200 and more computer time.

3. Some Critical Observations

Comparing, in some Montecarlo studies, different correlation coefficients when the underlying bivariate distribution is normal, the polychoric correlation is shown to be "the best in the sense of being closest to the true correlation" (Quiroga (1992)) [even if in a Montecarlo study Babakus et al. (1987) show that the polychoric coefficient provides the best estimates of model parameters and the worst fit statistics]. However, Muthen (1984), Aish and Jöreskog (1990) and Quiroga (1992) say also that "the assumption of underlying bivariate normality is too strong for most ordinary variables used in social sciences". There are not many studies on the use of polychoric coefficient with underlying not normal distribution. Lee and Lam (1988) study the robustness of polychoric coefficient only when the underlying distribution is elliptical (containing multivariate normal, platycentric and leptocentric distributions). Lee et al. (1995) even though obtaining quite satisfactory results about such robustness with moderate size random samples, say that "to draw a non definite conclusion, a longer simulation is needed". Moreover in literature only Quiroga (1992) tries, in a systematic way, to extend the distributive hypothesis of the continuous variables underlying the qualitative variables. First of all, she says that when you leave the normality assumptions, such variables are surely not consistent and slightly biased (Quiroga (1992)). Moreover, in a Montecarlo study she verifies that, when the underlying distribution is a skew-normal bivariate, the polychoric coefficient is the best choice only for sample size of 200-400 and large number of categories (5-9) and underestimates the true correlation coefficient. Then, by means of a measure of not normality, she verifies that when the underlying distributions are generated by the Fleishman-Vale-Maurelli polynomial transformation (with departure from normality due to skewness and kurtosis) the polychoric coefficient is robust, but overestimates the true correlation. Finally she proposes an extended polychoric coefficient with distribution given by a mixture of a normal and univariate skew-normal density function but she does not give empirical verifications. Therefore, until now, there are neither theoretical demonstrations nor empirical simulations which give satisfactory and generally valid reasons for using polychoric coefficient with underlying not normal distribution. Moreover the solutions based on polychoric coefficients: are not sensible to different scale of qualitative variables because they generally deal only with ordinal variables, reach solutions from variance covariance matrix of the reduced model different from the solutions obtained from the observed variables of the original measurement models and have the same problems of non identification of parameters and indeterminacy of latent variables of structural models with quantitative variables (Vittadini 1989).

4. An Alternative Proposal

In the quantitative case the problems of not uniqueness of the solutions are resolved by using linear combinations instead of causal latent variables of the observed variables (Wold (1982), Haagen and Vittadini (1991)). In this paper we propose to obtain the latent variables of the model as linear combinations of the observed mixed variables simultaneously quantified, by means of methods of multidimensional scaling using simultaneously optimal scaling and ordinary least squares method. In fact such methods resolve the problem of not normality of the variables because are distribution free (Young (1981)) and give unique solutions once chosen the method of multidimensional scaling. In order to avoid subjective choices about methodologies of multidimensional scaling (and therefore subjective solutions), we propose: to quantify the qualitative variable and to obtain the linear combinations of them by means of a unique objective function; to reach flexible solutions as regards to different kinds of linear combinations requested by the problems; to take into account the different scale of the qualitative variable. Therefore, among the variety of multidimensional scale methods we choose the family of Alsos method (Young (1981), and Keller and Wansbeek (1983)). These methods are based on alternating optimal scaling which quantifies qualitative variables and ordinary least squares which reach linear combinations of them in a iterative way. So they obtain solutions in a different way along the scale of the variables (ordinal-nominal, continuous-discrete), and the aim of analysis (e.g. Principal components, canonical correlation) giving answers to previous problems. Among the family of Alsos methods we avoid methods such as OSMOD (Saito and Otsu (1988)) or INDOMIX-CAMIX (Kiers (1991)) which obtain solutions in two stages. Instead we choose methods that simultaneously obtain quantifications of qualitative variables and their linear combinations (ADDALS (De Leeuw, Young, Takane (1976)) MORALS CORALS (Young, De Leeuw Takane, (1976)), PRINCALS (De Leeuw and Van Rijckevorsel (1980)) OVERALS (Van Der Burg and De Leeuw (1988)) respectively from the perspective of variance analysis, canonical correlation, principal components, multiple correspondence analysis. Moreover in order to obtain the latent variables from their real indicators as in Wold (1982), we apply the chosen Alsos methods to the subsets of mixed variables Y_β, X_δ characterized by submatrices $l_{y(\beta,\cdot)}, l_{x(\delta,\cdot)}$ with coefficients all different from zero. So we simultaneously obtain the quantification Y_β^*, X_δ^* of such mixed variables Y_β, X_δ and their linear transformations $\tilde{\eta}_\beta, \tilde{\xi}_\delta$ according to different aim of the analysis (e.g. canonical correlation, principal analysis etc.). Then in order to take into account the restrictions:

$$\begin{aligned} \text{cov}(\eta_\beta, \eta_\pi) &= 0; \text{cov}(\xi_\delta, \xi_\gamma) = 0; \text{cov}(\delta_\beta, \delta_\gamma); \text{cov}(u_{\delta_1}, u_{\delta_2}) = 0; \\ b_{(\beta, \mu)} &= 0; \gamma_{(\delta, \varphi)} = 0; l_{y(\beta, \alpha)} = 0; l_{x(\delta, \nu)} = 0; \end{aligned} \quad (2)$$

we obtain by means of the Restricted Regression Component Decomposition (RRCD) of quantified variables Y^*, X^* (Haagen and Vittadini (1998)) by means of an iterative process:

$$\begin{aligned} \eta_\beta^+ &= Q_{\eta_\mu \cup \eta_\pi} \eta_\beta^0 \quad (\eta_\beta^0 = Q_{\tilde{\eta}_\pi} \tilde{\eta}_\beta \quad (\beta \neq \pi)); * \eta_\mu = Q_{H_{(\beta, \mu)}} \eta_\mu^0; * y_\alpha = Q_{H_{(\beta, \alpha)}} y_\alpha \\ \xi_\delta^+ &= Q_{\eta_\gamma \cup \xi_\gamma} \xi_\delta^0 \quad (\xi_\delta^0 = Q_{\tilde{\xi}_\gamma} \tilde{\xi}_\delta \quad (\gamma \neq \delta)); * \eta_\gamma = Q_{H_{(\gamma, \delta)} \cup \Xi_{(\delta)}} \eta_\gamma^0; * x_\nu = Q_{\Xi_{(\delta)}} x_\nu \\ \tilde{\varepsilon}_j &= Q_{\tilde{H}_{(j)} \cup \tilde{\Xi}} \tilde{\eta}_j \quad (j=1, \dots, j-1); \delta_{\beta_1}^0 = Q_{(H^0 \cup \Xi^0 \cup \gamma_{\beta_1})} y_{\beta_1}; u_{\delta_1}^0 = Q_{(X \cup \Xi^0 \cup x_{\delta_1})} x_{\delta_1} \end{aligned} \quad (3)$$

where $H_{(\mu, \beta)}^0$ are the H^0 without η_μ^0, η_β^0 . $Q_{\tilde{\eta}_\pi}$ is the complement orthogonal to the orthogonal projector on the space generated by $\tilde{\eta}_\pi$ and the other symbols are defined in a similar way. So we have the following RRCD of $y_{\beta_1}, x_{\delta_1}, \eta_\beta^0$:

$$\begin{aligned} Q_{\Xi^0 \cup \gamma_{\beta_1}^0 / H^0} y_{\beta_1} &= P_{H^0} y_{\beta_1} + Q_{H^0 \cup \Xi^0 \cup \gamma_{\beta_1}} y_{\beta_1}; Q_{\gamma_{\delta_1} / \Xi^0} x_{\delta_1} = P_{\Xi^0} x_{\delta_1} + Q_{\gamma_{\delta_1} \cup \Xi^0} x_{\delta_1}; \\ \eta_\beta^0 &= P_{H_{(\beta)}} \eta_\beta^0 + P_{\Xi^0 / H_{(\beta)}} \eta_\beta^0 + Q_{H_{(\beta)} \cup \Xi^0} \eta_\beta^0 \end{aligned} \quad (4)$$

5. Numerical Example

The following variables are observed on a sample of 150 families casually chosen from 4103 american families that have been codified in the Federal Reserve Board research regarding National Income and Wealth of 1983.

Y_1 Job contract household (y_{11}), spouse (y_{12}); Occupation kind household (y_{13}), spouse (y_{14}); Occupation sector household (y_{15}), spouse (y_{16}). Y_2 Total health (y_{21}); Income (y_{22}); Debt (y_{23}). X_1 Age household (x_{11}), spouse (x_{12}); Sex household (x_{13}), spouse (x_{14}); Number of children (x_{15}); Race (x_{16}); Residence region (x_{17}); Civil Status (x_{18}); X_2 Educational Level household (x_{21}), spouse (x_{22}); Full time job years household (x_{23}), spouse (x_{24}); Part time job years household (x_{25}), spouse (x_{26}); Latent variables: labour force (η_1); Health and income (η_2), Civil status (ξ_1), Instruction grade (ξ_2).

In order to verify the causal dependence of the latent variables H from the latent variables Ξ we use the alternative proposal shown in paragraph 4 with the following restrictions: $l_{x_{(1,2)}, \eta_2} = 0, l_{y_{(2,1)}, \eta_2} = 0, l_{x_{(1,2)}, \xi_1} = 0, l_{x_{(2,1)}, \xi_2} = 0, \text{cov}(\Delta_1, \Delta_2) = 0,$

$\text{cov}(U_1, U_2) = 0$. The qualitative variables are quantified and the latent variables

are obtained as principal components with Princals method, the restrictions are then taken into account by means of RRCD. The variance-covariance of the observed variables and the results are shown in table 1.

Table 1: *The alternative proposal for the sample of American families.*

S_T									
0.0293	-0.0284	-0.0085	0.0248	-0.0047	0.0580	19.853	-0.0165	-0.0783	
-0.0284	8.2791	-1.2353	0.0917	0.8073	-0.3585	-79.107	12.671	0.5563	
-0.0085	-1.2353	17.0924	-0.0113	-0.6148	4.1994	-325.60	-17.426	-1.1561	
0.0248	0.0917	-0.0113	0.1420	0.0107	-0.0858	80.249	8.5699	-0.0821	
-0.0047	0.3073	-0.6148	0.0107	2.5357	-0.7303	-162.35	15.437	0.4694	
0.0580	-0.3585	4.1994	-0.0858	-0.7303	11.735	-65.533	-37.510	0.1070	
19.853	-79.107	-325.606	80.249	-0.0162	-65.533	320921	42682	-377.51	
-0.0165	12.671	-17.426	8.5699	15.437	-37.510	42622.1	17114	-21.173	
-0.0783	0.5563	-1.1561	-0.0821	0.4574	0.1070	-377.51	-21.173	7.4133	
S_Z									
1603.4	-16.291	59262	12.573	-16.997	-30333	37.233	-11617	9.074	-3.1692
-16.291	169.04	-53373	35773	07344	-11737	7.0739	03466	-23357	-2.8791
0.5926	-0.5337	9.3333	-5.1370	-0.0330	-1.1111	1.2672	-2.2332	-9.2626	-0.1623
12.573	8.0773	-13769	32374	-0.1312	-7.0707	-10200	-2.2376	-3.5646	-0.3074
-16.997	07344	-33037	-5.1312	11162	-57575	-2.3432	0.3512	2.5414	0.0965
-0.3033	-0.1173	-1.1111	-5.0607	-0.0637	1.0660	0.3026	-0.0637	1.5434	0.0169
37.233	7.0739	1.2672	-1.0200	-2.3432	3.0161	197.33	-1.3076	-2.5217	-8.4634
-11.617	03466	-3.3332	-2.2376	0.3512	-3.2393	-1.3076	0.2453	-3.3333	0.0719
0.9047	-2.5367	-9.2626	-0.3664	0.2041	1.5454	-2.5217	-0.0637	4.0091	2.4420
-3.1692	-2.8791	-0.1623	0.3074	0.0965	1.0660	-3.4634	0.0719	2.4420	4.2763
1.8004	107.29	-4.5161	7.7005	-0.2343	-17.121	10.037	0.4641	-57.035	-4.3171
-2.0307	-2.1334	2.4767	-0.9167	0.6135	-1.6969	4.1361	0.1967	1.2327	0.2771
34.131	45014	6.6269	3.3415	-0.7796	-37070	14.046	-0.2341	-6.7161	0.4076
-69.770	14.060	-9.7171	0.3945	1.2531	-5.4747	0.2234	0.3446	-3.4020	0.246
S_{D_1}									
0.0375	-0.0033	-0.0707	0.0207	0.0137	0.0483				
-0.0033	4.5237	1.1799	-0.0353	-0.2292	2.9146				
-0.0707	1.1799	11.0353	0.0137	0.1932	0.4736				
0.0207	-0.0353	0.0137	0.1326	0.0123	0.0542				
0.0137	-0.2292	0.8932	0.0123	2.1494	0.1236				
0.0483	2.9146	0.4736	0.0542	0.1236	10.7673				
S_{D_2}									
143417.7	6456.7	118.6							
6456.7	7047.8	81.98							
118.6	81.93	3.442							
S_{U_1}									
4507216	-2.2195	-0.4063	0.9215	-2.2092	-0.4451	-147014	-1.1134		
-2.2195	95.5067	-0.4451	7.6963	-0.2513	-0.2854	9.7268	0.3402		
-0.4063	-0.4451	0.2841	-0.1741	0.2133	-0.2256	0.7438	-0.0234		
0.9215	7.6963	-0.1741	2.7813	0.254	-0.2132	-5.5128	-0.0232		
-2.2092	-0.2513	0.2133	0.254	0.2256	-0.2256	-0.254	0.0232		
-0.4451	-0.2854	-0.0236	-0.0132	-0.0232	0.0232	0.0232	0.0232		
-147014	9.7268	0.7438	-5.5128	-0.254	0.0232	12.5576	0.0117		
-1.1134	0.3402	-0.0234	-0.0232	0.0232	0.0232	0.0117	0.0232		
S_{U_2}									
1.0289	0.1470	0.5453	-1.0289	1.7253	-0.9136				
0.1470	1.3709	1.3389	-1.5341	3.7266	-0.5184				
0.5453	1.3389	10.5882	-2.0286	1.2425	0.2277				
-1.0289	-1.5341	2.0286	10.4653	2.974	-1.0289				
1.7253	3.7266	1.2425	2.974	4.5834	-10.1886				
-0.9136	-0.5184	0.2277	-1.0289	-10.1886	23.5734				

With the example we can verify that the alternative proposal respects all the properties of the structural model described in paragraph 1 and the restrictions indicated in this paragraph. But the alternative proposal obtains unique

solutions solving all the problems of non-identification of parameters and indeterminacy of latent variables of the structural models with qualitative variables.

References

- Aish, A. M., Jöreskog, K. G. (1990). A panel model for political efficacy and responsiveness: an application of LISREL7 with weighted least squares, *Quality and Quantity*, 24, 405-426.
- Babakus, E., Ferguson, C.E. Jr, Jöreskog K.G. (1987). The sensitivity of confirmatory maximum likelihood factor analysis to violations of measurement scale and distributional assumptions, *Journal of Marketing Research*, 24, 222-228.
- De Leeuw, J., Van Rijckevorsel, J.(1980). *Homals & princals, some generalizations of principal components analysis*, In E. Diday et al. (Eds.), *Data Analysis and Informatics*, North-Holland Publishing Company, 231-241.
- De Leeuw, J., Young, F.W., Takane Y. (1976). Additive structure in qualitative data: an alternating least squares method with optimal scaling features, *Psychometrika*, 41, 471-503.
- Haagen, K., Vittadini, G. (1991). Regression Component Decomposition in Structural Analysis, *Communications in Statistics*, 20, 1153-1161.
- Haagen, K., Vittadini, G. (1998). Regression Component Decomposition Restricted. Un'alternativa al Lisrel model, *Metron*, 56, 1-2, in corso di pubblicazione.
- Jöreskog, K.G. (1990). New developments in Lisrel: analysis of ordinal variables using polychoric correlations and weighted least squares, *Quality and Quantity*, 24, 387-404.
- Jöreskog, K.G. (1994). On the estimation of polychoric correlations and their asymptotic covariance matrix, *Psychometrika*, 59, 3, 381-389.
- Kiers, H. A. L. (1991). Simple structure in component analysis techniques for mixtures of qualitative and quantitative variables, *Psychometrika* 56, 2, 197-212.
- Keller, W. J., Wansbeek, T. (1983). Multivariate methods for quantitative and qualitative data, *Journal of Econometrics*, 22, 91-111.
- Lee, S.Y, Lam, M. L.(1988). Estimation of polychoric correlation with elliptical latent variables. *Journal of statistic Computation and Simulation*, 30, 173-188.
- Lee, S.Y., Poon, W.Y., Bentler, P.M. (1990). Full maximum likelihood analysis of structural equation models with polytomous variables, *Statistics and Probability Letters*, 9, 91-97.

- Lee, S.Y., Poon, W.Y., Bentler, P.M. (1995). A two-stage estimation of structural equation models with continuous and polytomous variables, *British Journal of Mathematical and Statistical Psychology*, 48, 339-358.
- Muthen, B. (1984). A general structural equation model with dichotomous, ordered categorical, and continuous latent variable indicators, *Psychometrika*, 49, 1, 115-132.
- Quiroga, A.M. (1992). *Studies of the Polychoric Correlation and other Correlation Measures for Ordinal Variables*, PhD thesis, Uppsala University.
- Rigdon, E.E., Ferguson, C.E. Jr. (1991) The performance of the polychoric correlation coefficient and selected fitting functions in confirmatory factor analysis with ordinal data, *Journal of Marketing Research*, 28, 491-497.
- Saito, T., Otsu, T. (1988). A method of optimal scaling for multivariate ordinal data and its extensions, *Psychometrika*, 53, 1, 5-25.
- Van Der Burg, E., De Leeuw, J.(1988). Homogeneity analysis with k sets of variables: an alternating least squares method with optimal scaling features, *Psychometrika*, 53,2, 177-197.
- Vittadini, G. (1989). *Indeterminacy Problems in the LISREL Model*, in *Multivariate Behavioral Research*, Forth Worth (Texas), 24, 4, 397-414.
- Wold, H. (1982). *Soft Modelling: the basic design and some extensions*, in Jöreskog K.G., Wold H., *Systems under indirect observation: casuality, structure, prediction*, North - Holland, Amsterdam, 2, 1-54.
- Young, F.W. (1981). Quantitative analysis of qualitative data, *Psychometrika*, 46, 357-388.
- Young, F.W., De Leeuw, J., Takane, Y. (1976). Regression with qualitative and quantitative variables: an alternating least squares method with optimal scaling features, *Psychometrika*, 41, 4, 505-529.