



Covariates and ordinal data

Contents

1.	INTRODUCTION	1
2.	UNIVARIATE PROBIT REGRESSION	2
3.	UNIVARIATE LOGIT REGRESSION	7
4.	TESTING THE MODEL	8
5.	BIVARIATE PROBIT REGRESSION	9
6.	MULTIVARIATE PROBIT REGRESSION	10
7.	PRELIS IMPLEMENTATION	12
8.	A SMALL EXAMPLE	13
9.	DATA SCREENING	14
10.	PROBIT REGRESSION OF NOSAY	15
11.	PROBIT AND LOGIT REGRESSION OF ALL EFFICACY VARIABLES	19
12.	ESTIMATING THE JOINT COVARIANCE MATRIX	22
13.	A MIMIC MODEL FOR EFFICACY AND RESPONSES	24

1. Introduction

This example is the fifth of a set of examples extracted from a note by K.G. Jöreskog first posted on the SSI website in 2005 with the title “*Structural Equation Modeling with Ordinal Variables using LISREL*”.

In the previous four examples I assumed that all observed variables were ordinal. Thus, in the second example, I described the analysis of ordinal variables in cross-sectional studies, in the third, I described the analysis of ordinal variables in longitudinal studies, and in the fourth, I described the analysis of ordinal variables observed in several groups. In this example I consider the case when one or more ordinal variables are observed jointly with a set of possibly explanatory variables, so called covariates. These covariates can be dummy-coded categorical variables or measured variables on an interval scale. They are assumed not to contain measurement error. With PRELIS one can estimate the effect of the covariates on the probability of response in various categories of the ordinal variables using either the probit or the logit model. PRELIS can also estimate the joint covariance matrix of the covariates and the variables underlying the ordinal variables. This can be used for further modeling in LISREL.

Continuing my analysis of the Efficacy variables from the Political Action Survey, I illustrate the analysis of the six Efficacy variables using four covariates: Gender, Age, Education, and a Left-Right Scale. For information about the Political Action Survey and the Efficacy variables, see the first example. As in the second and third examples, I will only use the data from the USA sample. Since probit and logit regression have not been well documented in the LISREL literature, I give a rather technical description in Sections 2-6. Readers who are merely interested in how to do it with LISREL can skip this and proceed to Section 7.

2. Univariate Probit Regression

Let y be a single ordinal variable with m categories and let $\mathbf{x}(q \times 1)$ be a vector of covariates. The term univariate is used here in the sense of one variable at a time. One can very well have several ordinal variables but they are analyzed one at a time. Corresponding to y there is an underlying continuous variable y^* . The connection between the ordinal variable y and the underlying variable y^* is

$$y = a \Leftrightarrow \tau_{a-1} < y^* < \tau_a, \quad a = 1, 2, \dots, m, \quad (17)$$

where $\tau_0 = -\infty, \tau_1 < \tau_2 < \dots < \tau_{m-1}, \tau_m = +\infty$, are threshold parameters. For variable y with m categories, there are $m - 1$ threshold parameters $\boldsymbol{\tau} = (\tau_1, \tau_2, \dots, \tau_{m-1})$.

The specification (17) is the same as in the second example. In fact, the development there is a special case of the more general case described here, namely when $q = 0$.

Consider the regression of y^* on \mathbf{x} :

$$y^* = \alpha + \boldsymbol{\gamma}' \mathbf{x} + z, \quad (18)$$

where α is an intercept term, $\boldsymbol{\gamma}$ is a vector of regression coefficients, and z is an error term. The underlying variable y^* is not observed; only the ordinal variables y and \mathbf{x} are.

The probit model assumes that z is normally distributed with mean 0 and variance ψ^2 , i.e., y^* is normal conditional on \mathbf{x} :

$$y^* \sim N(\alpha + \boldsymbol{\gamma}' \mathbf{x}, \psi^2).$$

It follows that the probability $\pi_c(\mathbf{x})$ of a response in category c or lower, conditional on \mathbf{x} , where $c = 1, 2, \dots, m-1$, is

$$\pi_c(\mathbf{x}) = \Phi\left(\frac{\tau_c - \alpha - \boldsymbol{\gamma}' \mathbf{x}}{\psi}\right), \quad (19)$$

where Φ is the standard normal distribution function.

Equation (19) can be viewed as a special case of a generalized linear model, see e.g., McCullagh & Nelder (1983). In this tradition there are no concepts of underlying variables and thresholds. Instead, (19) is written

$$\pi_c(\mathbf{x}) = \Phi(\alpha_c^* - \gamma' \mathbf{x}), \quad (20)$$

where

$$\alpha_c^* = \psi^{-1}(\tau_c - \alpha),$$

is interpreted as an intercept term and

$$\gamma^* = \psi^{-1} \gamma,$$

is a vector of regression coefficients.

To explain the term *probit regression*, take the inverse of (20):

$$\Phi^{-1}[\pi_c(\mathbf{x})] = \alpha_c^* - \gamma' \mathbf{x},$$

where Φ^{-1} is the inverse function of Φ . The quantity $\Phi^{-1}(\pi)$ is called the probit of π . If π goes from 0 to 1, $\Phi^{-1}(\pi)$ goes from $-\infty$ to $+\infty$. Equation (21) shows that the probit of $\pi_c(\mathbf{x})$ is linear in \mathbf{x} , hence the term probit regression. Note that the sign of γ is negative in (20) but positive in (18). In (20), for example, if γ_1 is positive, the probability of a response in category c or lower decreases as x_1 increases. Which says the same thing as the probability of a response in a category higher than c increases with x_1 . In (18), however, y^* increases if x_1 increases which increases the probability of a higher response. Thus, the two models are equivalent.

One can regard (21) as $m - 1$ parallel regression lines. Note that the intercepts vary with c but the regression coefficients are the same. For ordinal variables the intercepts must satisfy the order condition

$$\alpha_1^* < \alpha_2^* < \dots < \alpha_{m-1}^*.$$

To illustrate the function $\pi_c(\mathbf{x})$ in (19), consider the case of a single covariate x . Let $\alpha = 0$, $\psi = 1$ and denote

$$\pi_{\tau\gamma}(x) = \Phi(\tau - \gamma x).$$

Fig. 7 shows four curves $\pi_{\tau\gamma}(x)$ for $-10 < x < 10$ using the parameter values

- Curve 1** $\tau = -0.5$ and $\gamma = 1.0$
Curve 2 $\tau = 1.5$ and $\gamma = 1.0$
Curve 3 $\tau = -0.5$ and $\gamma = 0.4$
Curve 4 $\tau = 1.5$ and $\gamma = 0.4$.

It is seen that the probability of a response in category c or lower, i.e., the probability that $y^* \leq \tau$, decreases with x . The larger is, the faster is the rate of decrease. As τ increases or decreases, the curves are just shifted vertically.

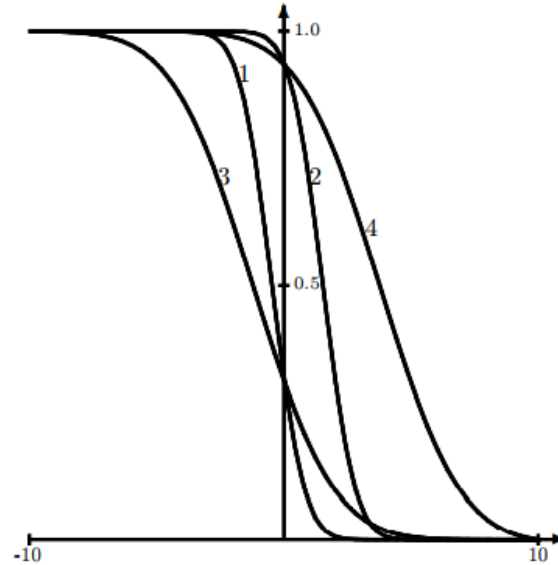


Figure 7: Four Cumulative Response Functions

I now return to the development of equation (19). There are two fundamental indeterminacies in (19):

- One can add a constant to all τ 's and to α .
- One can multiply all τ 's, α , and γ by a constant and multiply ψ by the same constant.

Neither of these changes has any effect on the right-hand side of (19). This is a reflection of the fact that since only ordinal information is available about y , y^* is only determined up to a linear transformation (Actually, y^* is only determined up to a monotonic transformation, but under normality the transformation must be linear.)

PRELIS has two ways of resolving these indeterminacies:

Standard Parameterization: $\alpha = 0$ and $\psi = 1$

Alternative Parameterization: $\tau_1 = 0$ and $\tau_2 = 1$

These parameterizations fix the origin and unit of measurement of y^* in two different ways. The Standard Parameterization is the same as used in Generalized Linear Models. The Alternative Parameterization requires that $m \geq 3$. If $m = 2$ under this parameterization, PRELIS will set $\tau_1 = 0$ and $\psi = 1$.

For $m \geq 3$, the parameters of the two parameterizations are given in the following table.

Parameterization	Intercept	Error Var.	Thresholds					Regr. Coeff.			
Standard	0	1	τ_1	τ_2	τ_3	\cdots	τ_{m-1}	γ_1	γ_2	\cdots	γ_q
Alternative	α	ψ^2	0	1	τ_3^*	\cdots	τ_{m-1}^*	γ_1^*	γ_2^*	\cdots	γ_q^*

where

$$\begin{aligned}\alpha &= -\tau_1 / (\tau_2 - \tau_1), \quad \psi = 1 / (\tau_2 - \tau_1), \\ \tau_i^* &= (\tau_i - \tau_1) / (\tau_2 - \tau_1), \quad i = 3, 4, \dots, m-1, \\ \gamma_i^* &= \gamma_i / (\tau_2 - \tau_1), \quad i = 1, 2, \dots, q.\end{aligned}$$

It should be emphasized that the two parameterizations are equivalent in the sense that there is a one-to-one correspondence between the two sets of parameters.

For estimation, the probability of a response in category a is needed, where $a = 1, 2, \dots, m$. This is

$$\Pr\{y = a \mid \mathbf{x}\} = \pi_a(\mathbf{x}) = \Phi\left(\frac{\tau_a - \alpha - \gamma' \mathbf{x}}{\psi}\right) - \Phi\left(\frac{\tau_{a-1} - \alpha - \gamma' \mathbf{x}}{\psi}\right) \quad (22)$$

It is convenient to refer to (22) as the category probability function and to (19) as the cumulative probability function.

For a single x , the category probability functions in (22) are shown in Fig. 8 for $\alpha = 0$, $\psi = 1$, $\tau_1 = -0.5$, $\tau_2 = 1.5$ and $\gamma = 1$ (Curve 1) and $\gamma = 0.4$ (Curve 2). As x increases the category probability increases up to a maximum

$$2\Phi\left[\frac{1}{2}(\tau_2 - \tau_1)\right] - 1 \quad \text{at} \quad x = \frac{\tau_1 + \tau_2}{2\gamma},$$

and then decreases. The rate of increase and decrease is larger for larger than for smaller γ . Note that the maximum is independent of γ .

Suppose we have a random sample of N independent observations of y and \mathbf{x} :

$$(y_i, \mathbf{x}_i), \quad i = 1, 2, \dots, N.$$

Let $k_{ia} = 1$, if $y_i = a$, and $k_{ia} = 0$, otherwise. Then

$$E(k_{ia} | \mathbf{x}_i) = \Phi\left(\frac{\tau_a - \alpha - \gamma' \mathbf{x}_i}{\psi}\right) - \Phi\left(\frac{\tau_{a-1} - \alpha - \gamma' \mathbf{x}_i}{\psi}\right) = \pi_{ia}^*(\mathbf{x}_i), \quad (23)$$

say. The likelihood of the sample is

$$L = \prod_{i=1}^N \left\{ \prod_{a=1}^m [\pi_{ia}^*(\mathbf{x}_i)^{k_{ia}}] \right\} p(\mathbf{x}_i), \quad (24)$$

where $p(\mathbf{x})$ is the density function of \mathbf{x} . The latter is unspecified and assumed to have no parameters of interest. The parameter vector is

$$\boldsymbol{\theta} = (\boldsymbol{\tau}, \alpha, \boldsymbol{\gamma}, \psi).$$

This can be estimated by maximizing the likelihood L of either of the two parameterizations.

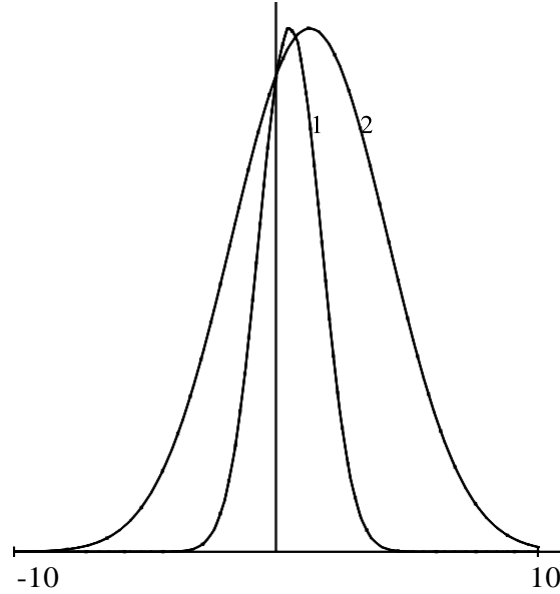


Figure 8: Two Category Response Functions (probit)

3. Univariate Logit Regression

One can obtain *logit regression*, sometimes called *logistic regression*, in the same way simply by replacing the normal distribution function $\Phi(x)$ by the logistic distribution function

$$\Psi(u) = \frac{e^u}{1 + e^u}.$$

The inverse function of Ψ is

$$\Psi^{-1}(\pi) = \ln \frac{\pi}{1 - \pi}.$$

The quantity $\ln \frac{\pi}{1 - \pi}$ is called the *logit* of π . If π goes from 0 to 1, $\text{logit}(\pi)$ goes from $-\infty$ to $+\infty$.

The logit model is

$$\ln \frac{\pi_c(\mathbf{x})}{1 - \pi_c(\mathbf{x})} = \alpha_c^* - \gamma^* \mathbf{x}. \quad (25)$$

This is also a special case of a Generalized Linear Model, see McCullagh & Nelder (1983). The logit model seems to be more often used in practice than the probit model. This is probably because Ψ^{-1} has an explicit form. However, with computers it is almost as easy to compute Φ^{-1} as it is to compute Ψ^{-1} .

PRELIS estimates the logit model in the form

$$\pi_c(\mathbf{x}) = \Psi\left(\frac{\tau_c - \alpha - \gamma' \mathbf{x}}{\psi}\right), \quad (26)$$

using either the Standard Parameterization or the Alternative Parameterization as defined in the previous Section.

The probability of a response in category a is

$$\Pr\{y = a \mid \mathbf{x}\} = \pi_a(\mathbf{x}) - \pi_{a-1}(\mathbf{x}) = \Psi\left(\frac{\tau_a - \alpha - \gamma' \mathbf{x}}{\psi}\right) - \Psi\left(\frac{\tau_{a-1} - \alpha - \gamma' \mathbf{x}}{\psi}\right), \quad (27)$$

This probability as a function of a single x is shown in Fig. 9 for the same parameters as in Fig. 8. It is seen that the logit model gives less probability to the category corresponding to $\tau_1 < y^* < \tau_2$ and more probability to the other categories than the probit model. In general, the logit model gives less probability to the middle categories and more probability to the outer categories than the probit model.

The logistic and the normal distribution are similar, but the variance of the logistic distribution $\Psi(u)$ is not 1 as is the case of the normal distribution $\Phi(u)$. Lord & Novick (1968, p. 299) noted that

$$|\Phi(u) - \Psi(1.7u)| < 0.001, \quad \forall x. \quad (28)$$

Because of this closeness, results obtained under the Standard Parameterization with the logit model are likely to be close to those obtained with the probit model except for a scale factor. As will be demonstrated, the Alternative Parameterization eliminates this scale factor and makes the regression equations directly comparable and similar.

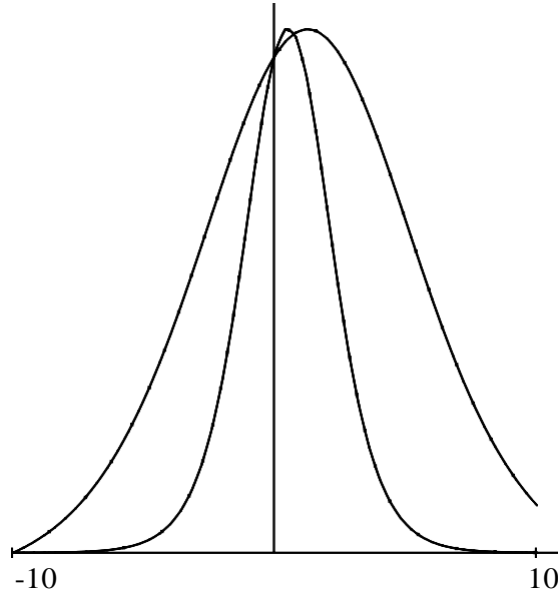


Figure 9: Two Category Response Functions (logit)

4. Testing the Model

The univariate probit and logit models make strong assumptions about the cumulative response functions in the form of (19) and (26), respectively. Can these assumptions be tested?

The typical way of doing this is to compute a deviance, i.e., the difference between $-2\ln \hat{L}$ for the model and the same quantity for another more general model, where \hat{L} is the maximum value of (24). This kind of deviance has not been implemented in PRELIS because it is not obvious what the more general model should be. However, PRELIS prints the value of $-2\ln \hat{L}$ so one can compare this for different models. We have also implemented a test of the hypothesis that $\gamma = \mathbf{0}$, i.e., that all regression coefficients are zero. This can also be regarded as a measure of how much better the model fits than the model with no covariates. This will be illustrated in Section 11.

5. Bivariate Probit Regression

The normal distribution generalizes naturally to the bivariate and multivariate case. The logistic distribution function, however, does not have any convenient generalization to the bivariate and multivariate case. For this reason I consider only the case of underlying bivariate and multivariate normality in what follows.

Consider two ordinal variables y_g and y_h with underlying continuous variables y_g^* and y_h^* , respectively. The equations to be estimated are

$$y_g^* = \alpha_g + \gamma_g' \mathbf{x} + z_g, \quad (29)$$

$$y_h^* = \alpha_h + \gamma_h' \mathbf{x} + z_h, \quad (30)$$

where α_g and α_h are intercept terms, γ_g and γ_h are vectors of regression coefficients, and z_g and z_h are error terms. It is assumed that z_g and z_h have a bivariate normal distribution with means zero and covariance matrix

$$\begin{pmatrix} \psi_g^2 & \psi_{gh} \\ \psi_{gh} & \psi_h^2 \end{pmatrix}.$$

In the Standard Parameterization this is a correlation matrix with correlation ρ_{gh} . Variable y_g has thresholds

$$\tau_g = (\tau_{g,1}, \tau_{g,2}, \dots, \tau_{g,m_g-1}),$$

and variable y_h has thresholds

$$(\tau_{h,1}, \tau_{h,2}, \dots, \tau_{h,m_h-1}).$$

The probability that an individual i with covariates \mathbf{x}_i responds in category a on y_g and in category b on y_h is

$$\pi_{igh,ab} = \Pr\{y_{ig} = a, y_{ih} = b \mid \mathbf{x}_i\} = \int_{\tau_{ig,a-1}^*}^{\tau_{ig,a}^*} \int_{\tau_{ih,b-1}^*}^{\tau_{ih,b}^*} \phi^{(2)}(u, v, \rho_{gh}) du dv, \quad (31)$$

where

$$\tau_{ig,a}^* = \frac{\tau_{ig,a} - \alpha_g - \gamma_g' \mathbf{x}_i}{\psi_g} \quad (32)$$

and $\phi^{(2)}(u, v, \rho)$ is the density function of the standardized bivariate normal distribution with correlation ρ . The parameter vector is

$$\boldsymbol{\theta} = (\boldsymbol{\theta}_g, \boldsymbol{\theta}_h, \rho_{gh}),$$

where

$$\boldsymbol{\theta}_g = (\boldsymbol{\tau}_g, \alpha_g, \boldsymbol{\gamma}_g, \psi_g),$$

$$\boldsymbol{\theta}_h = (\boldsymbol{\tau}_h, \alpha_h, \boldsymbol{\gamma}_h, \psi_h).$$

The likelihood function is
$$L = \prod_{i=1}^N \left(\prod_{a=1}^{m_g} \prod_{b=1}^{m_h} \pi_{igh,ab}^{k_{igh,ab}} \right) p(\mathbf{x}_i), \quad (33)$$

where $k_{igh,ab} = 1$ if case i responds in category a on y_g and in category b on y_h , and $k_{igh,ab} = 0$, otherwise.

PRELIS estimates $\boldsymbol{\theta}_g$ and $\boldsymbol{\theta}_h$ from the univariate marginal distribution of y_g and y_h , respectively, as described in The first example in this set. Given these estimates, PRELIS estimates ρ_{gh} by maximizing the bivariate likelihood L in (33). Under the Alternative Parameterization, the conditional covariance between y_g^* and y_h^* is estimated as

$$\hat{\psi}_{gh} = \hat{\psi}_g \hat{\psi}_h \hat{\rho}_{gh}.$$

6. Multivariate Probit Regression

Let $\mathbf{y}(p \times 1)$ be a vector of ordinal variables with underlying variables \mathbf{y}^* . It is assumed that

$$\mathbf{y}^* | \mathbf{x} \sim N(\boldsymbol{\alpha} + \boldsymbol{\Gamma}\mathbf{x}, \boldsymbol{\Psi}).$$

The rows of $\boldsymbol{\alpha}$ and $\boldsymbol{\Gamma}$ and the diagonal elements of $\boldsymbol{\Psi}$ are estimated from the univariate margins as described in Section 1, and the off-diagonal elements of $\boldsymbol{\Psi}$ are estimated from the bivariate margins as described in Section 5.

Denoting these estimates as $\hat{\boldsymbol{\alpha}}$, $\hat{\boldsymbol{\Gamma}}$, and $\hat{\boldsymbol{\Psi}}$, we have the following:

- The estimated conditional covariance matrix of \mathbf{y}^* for given \mathbf{x} is $\hat{\boldsymbol{\Psi}}$. In the Standard Parameterization this is a correlation matrix.
- The estimated unconditional covariance matrix of \mathbf{y}^* is

$$\hat{\boldsymbol{\Gamma}} \mathbf{S}_{xx} \hat{\boldsymbol{\Gamma}}' + \hat{\boldsymbol{\Psi}},$$

where \mathbf{S}_{xx} is the sample covariance matrix of \mathbf{x} .

- The estimated joint unconditional covariance matrix of \mathbf{y}^* and \mathbf{x} is

$$\hat{\Sigma} = \begin{pmatrix} \hat{\Gamma} \mathbf{S}_{xx} \hat{\Gamma}' + \hat{\Psi} & \\ & \mathbf{S}_{xx} \end{pmatrix}. \quad (34)$$

The relationship between the Standard and Alternative Parameterizations can be expressed in matrix form as follows. Let \mathbf{D} be the diagonal matrix of order $p \times p$

$$\mathbf{D} = \text{diag} \left(\frac{1}{\tau_{1,2} - \tau_{1,1}}, \frac{1}{\tau_{2,2} - \tau_{2,1}}, \dots, \frac{1}{\tau_{p,2} - \tau_{p,1}} \right), \quad (35)$$

and let y_s^* and y_A^* denote the vector of underlying variables in the Standard and Alternative Parameterizations, respectively. Then

$$\hat{\Gamma}_A = \mathbf{D} \hat{\Gamma}_S, \quad (37)$$

$$\hat{\Psi}_A = \mathbf{D} \hat{\Psi}_S \mathbf{D}. \quad (38)$$

Using the same notation for the matrix $\hat{\Sigma}$ in (34), we have

$$\hat{\Sigma}_A = \mathbf{D}_1 \hat{\Sigma}_S \mathbf{D}_1, \quad (39)$$

where \mathbf{D}_1 is the diagonal matrix of order $p + q \times p + q$

$$\mathbf{D}_1 = \text{diag} \left(\frac{1}{\tau_{1,2} - \tau_{1,1}}, \frac{1}{\tau_{2,2} - \tau_{2,1}}, \dots, \frac{1}{\tau_{p,2} - \tau_{p,1}}, 1, 1, \dots, 1 \right). \quad (40)$$

PRELIS can also estimate the asymptotic covariance matrix of $\hat{\Sigma}$.

There is no latent variable model (LISREL model) imposed on the $\hat{\Sigma}$ in (34). It is an unconstrained covariance matrix just as a sample covariance matrix \mathbf{S} for continuous variables. It can therefore be used for modeling in LISREL just as if \mathbf{y}^* and \mathbf{x} were directly observed. This is illustrated in Section 14.

7. PRELIS Implementation

I illustrate the case of 3 ordinal variables and 4 covariates. Let Y1 Y2 Y3 be the names of the ordinal variables and let X1 X2 X3 X4 be the names of the covariates.

Probit regression of Y1 is obtained by the PRELIS command

```
PR Y1 on X1 X2 X3 X4
```

Similarly, logit regression of Y1 is obtained by the PRELIS command

```
LR Y1 on X1 X2 X3 X4
```

One can select any subset of y variables and any subset of x variables to be included in the equation. Thus, one can obtain the univariate probit or logit regression for all the ordinal variables simultaneously. For example,

```
PR Y1 Y2 Y3 on X1 X2 X3 X4
```

will give three univariate probit regressions. Note the word on (or ON) separating the ordinal variables from the covariates.

One can have several PR and/or LR commands in the same input file. All x -variables used as covariates must be declared continuous before the first PR or LR command, or else they must have at least 16 different values.

The Standard Parameterization is used by default. To obtain the Alternative Parameterization put AP on the Output line. The PR or LR command produces only univariate probit or logit regressions. Thus an MA value specified on the Output line has no meaning. To obtain the matrix $\hat{\Sigma}$ in (34), use an FI command and put MA = CM on the Output line. No other value of MA is meaningful since $\hat{\Sigma}$ is a covariance matrix even in the Standard Parameterization. There are two reasons why the covariance matrix $\hat{\Sigma}$ in (34) is not computable with PR or LR commands:

- Since one can have several PR or LR commands in the same PRELIS command file, there is no way PRELIS will know which covariance matrix to compute.
- Since the logistic distribution does not generalize to the multivariate case, the covariance matrix can only be estimated under multivariate normality. It would be odd to estimate the univariate parameters and under the logistic distribution and then estimate the covariance of the error terms under multivariate normality.

The various alternatives are illustrated in the sections that follow.

8. A Small Example

Before proceeding to analyze the Efficacy variables, I illustrate the various alternatives by means of a small example based on generated data. File **ORDATA.RAW** contains data in free format on one ordinal variable y and two covariates x_1 and x_2 . To estimate the probit regression in the Standard Parameterization, use the following PRELIS command file (file **ORDATA.PRL**):

```
Data Ninputvars = 3
Labels
Y X1 X2
Rawdata = ORDATA.RAW
Continuous X1 X2
LR Y on X1 X2
Output
```

The probit regression is estimated as

```
Thresholds:  -2.034-1.0870.537 1.925

Y = 1.006*X1 + 2.028*X2 + Error, R2 = 0.838 (0.0860) (0.119)
11.696 16.997
```

To estimate the same regression in the Alternative Parameterization, just put AP on the Output line. This gives the following results:

```
Thresholds:  0.0 1.0 2.773 4.251

Y = 2.147 + 1.061*X1 + 2.141*X2 + Error, R2 = 0.838 (0.0907)
(0.126)
11.696 16.997
```

Note that

- The regression coefficients in the Standard and Alternative Parameterizations are different but the t -values are the same.
- R^2 is the same.
- Although different, the regression coefficients are rather close. However, this is just a coincidence that occurs because $\hat{\tau}_2 - \hat{\tau}_1$ is close to 1.

To use logit regression, put LR instead of PR. Logit regression gives the following results in the Standard Parameterization:

```
Thresholds:  -3.639-1.9680.993 3.464

Y = 0.0 + 1.790*X1 + 3.631*X2 + Error, R2 = 0.943
Standerr      (0.154)      (0.227)
```

Comparing the standard solutions for probit and logit regression, it is seen that the regression coefficients are quite different. However, a closer look shows that the regression coefficients of the logit equation are approximately 1.79 times those of the probit regression. This confirms

the statement made earlier that the regression coefficients will be roughly proportional. The scale factor 1.79 may require some further explanation. The factor 1.7 in (28) should be regarded as an approximate population quantity, whereas the scale factor 1.79 is estimated from a random sample of 400 observations.

That the results of the probit and logit regressions are close can be seen much better if one uses the Alternative Parameterization. The result of logit regression in the Alternative Parameterization is:

Thresholds: 0.0 1.0 2.773 4.251

$$Y = 2.178 + 1.071 \cdot X_1 + 2.173 \cdot X_2 + \text{Error}, R^2 = 0.943$$

Standerr (0.0924) (0.136)

As can be seen, this is quite similar to the corresponding regression equation for the probit model.

9. Data Screening

I now return to the analysis of the *Efficacy* variables in the Political Action Survey described in the first example in this set. The data file for this illustration is **USA.RAW**. This contains 10 variables in free format. The first six are the six *Efficacy* variables; the other four variables are (the original variable names are given in parenthesis):

- YOB Year of birth with *Don't Know* coded as 1998 and *No Answer* coded as 1999 (V0146). Recall that the interviews were done in 1974.
- GENDER Gender coded as 1 for Male, 2 for Female, and 9 for *No Answer* (V0283).
- LEFTRIGH A left-right scale from 1 to 10 with *Don't Know* coded as 98 and *No Answer* coded as 99 (V0020).
- EDUCAT Education coded as 1 for Compulsory level only, 2 for Middle level, 3 for Higher or Academic level, and 9 for *No Answer* (V0214)

As always, it is a good idea to begin with a data screening. This can be done by running the following PRELIS command file (file **ORD51.PRL**)

```
Screening the Data in USA.RAW
Data Ninputvars = 10
Labels
NOSAY VOTING COMPLEX NOCARE TOUCH INTEREST YOB GENDER LEFTRIGH EDUCAT
Rawdata = USA.RAW
Clabls NOSAY - INTEREST 1=AS 2=A 3=D 4=DS 8=DK 9=NA
Clabls GENDER 1=MALE 2=FEMA 9=NA
Clabls LEFTRIGH 98=DK 99=NA
Clabls EDUCAT 1=COMP 2=MIDD 3=HIGH 9=NA
Output
```

The output reveals that

- There are 1719 cases in the USA sample, 736 males and 983 females.

- The marginal distributions of the six efficacy variables are those reported in the first example in this set.
- There are more than 15 different birthyears in the sample. The oldest person was born in 1882. Only one person did not report his/her birthyear and nobody reported not knowing his/her year of birth.
- As many as 547 persons or 31.8% did not place themselves on the left-right scale.
- Only 8 persons did not answer the education question.

Before one can proceed one must decide how to treat the *Don't Know* and *No Answer* responses. In The first example in this set, I discussed various alternative ways of dealing with missing values. I do not want to repeat that discussion here. The major difficulty is to decide how to treat the 547 people who did not answer the LEFTRIGH variable. Does this mean that these people are in the middle of the scale, or that the concept of left-right has no meaning for them, or what? I do not know. So I will treat them as having provided no information.

File **ORD51A.PRL** eliminates all cases with *Don't Know* and *No Answer* responses (listwise deletion) and saves the data on all complete cases in a PRELIS system file called **USA.PSF**. In addition, AGE is computed as 1974 - YOB. This is a proxy for age. The output file shows that the resulting listwise sample size is 1076. Thus, 643 cases were lost.

10. Probit Regression of NOSAY

The following PRELIS command file (file **ORD52.PRL**) estimates the probit regression of NOSAY (as a y-variable) on GENDER, LEFTRIGH, EDUCAT, and AGE (as x-variables - covariates) using the Alternative Parameterization.

```
Probit Regression of NOSAY
SY=USA.LSF
Continuous GENDER - AGE
PR NOSAY on GENDER - AGE
Output AP
```

The output gives the following information about the probit regression:

```
Univariate Probit Regression for    NOSAY
Alternative Parameterization

Thresholds: 0.0 1.0 2.825

      NOSAY = 0.173 + 0.00927*GENDER + 0.0437*LEFTRIGH + 0.371*EDUCAT
Standerr      (0.0696)          (0.0189)          (0.0534)
Z-values      0.133            2.308            6.946
P-values      0.894            0.021            0.000

      + 0.00467*AGE + Error, R2 = 0.059
      (0.00213)
      2.187
      0.029
```

Because the t -values for LEFTRIGH, EDUCAT, and AGE are all positive and larger than 2, this means that people on the right of the left-right scale, people with higher education and older people have a tendency to respond higher on the ordinal scale for NOSAY, that is, they are likely to disagree or disagree strongly to NOSAY. This seems quite plausible.

The corresponding logit regression is obtained by replacing PR with LR. The resulting regression equation is

```
Univariate Logit Regression for      NOSAY
Alternative Parameterization

Thresholds: 0.0 1.0 2.771

      NOSAY = 0.114 - 0.00285*GENDER + 0.0492*LEFTRIGH + 0.379*EDUCAT
Standerr      (0.0693)      (0.0188)      (0.0538)
Z-values      -0.0411      2.616      7.043
P-values      0.967      0.009      0.000

      + 0.00456*AGE + Error, R2 = 0.171
      (0.00212)
      2.147
      0.032
```

The logit regression is very similar to the probit regression, but note that R^2 is larger for the logit model than for the probit model. I will discuss the issue of the fit of the probit vs logit model in the next Section.

In addition to the ordinary output file **ORD52.OUT**, the run of **ORD52.PRL** gives information about the fit of the probit and logit regressions. For the initial probit regression of NOSAY, the result looks like this:

```
-2lnL for Full Model      2344.591
-2lnL for Intercept-Only Model      2398.103
Chi-Square for Testing Intercept-Only Model      53.512
Degrees of Freedom      4
```

This does not give much information; only that the four covariates fit much better than no covariate at all. However, consider entering the covariates stepwise one at a time using the following input file (**ORD52A.PRL**):

```
Probit Regression of NOSAY
SY=USA.LSF
Continuous GENDER - AGE
PR NOSAY on GENDER
PR NOSAY on GENDER AGE
PR NOSAY on GENDER AGE LEFTRIGH
PR NOSAY on GENDER AGE LEFTRIGH EDUCAT
Output AP
```

The output file now provides similar information for each regression

-2lnL for Full Model	2398.100
PR NOSAY on GENDER:	
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	0.004
Degrees of Freedom	1
PR NOSAY on GENDER AGE:	
-2lnL for Full Model	2395.940
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	2.163
Degrees of Freedom	2
PR NOSAY on GENDER AGE LEFTRIGH:	
-2lnL for Full Model	2393.094
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	5.009
Degrees of Freedom	3
PR NOSAY on GENDER AGE LEFTRIGH EDUCAT:	
-2lnL for Full Model	2344.591
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	53.512
Degrees of Freedom	4

This can be interpreted as follows. GENDER is no better than no covariate at all, i.e., GENDER alone cannot be used to predict NOSAY. If AGE is used together with GENDER there is no significant improvement in fit. GENDER and AGE alone does not predict NOSAY. If LEFTRIGH is added to the equation, there is still no significant improvement in fit because $5.009 - 2.163 = 2.846$ is not significant as a chi-square with one degree of freedom. If EDUCAT is added to the equation, there is a highly significant improvement in fit. This suggest that EDUCAT is the best predictor of NOSAY. These findings are confirmed in the output file **ORD52A.OUT**.

How come that LEFTRIGH and AGE are significant in the last equation whereas they are not significant in any equation that does not include EDUCAT? The reason is that EDUCAT is correlated with LEFTRIGH and AGE thereby generating interactive effects of LEFTRIGH and AGE. This can be seen by entering the covariates one at a time in the opposite order, see file **ORD52B.PRL**. The output from this run is

PR NOSAY on EDUCAT:	
-2lnL for Full Model	2357.968
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	40.136
Degrees of Freedom	1
PR NOSAY on EDUCAT LEFTRIGH:	
-2lnL for Full Model	2349.465
-2lnL for Intercept-Only Model	2398.103

Chi-Square for Testing Intercept-Only Model 48.638
 Degrees of Freedom 2

PR NOSAY on EDUCAT LEFTRIGH AGE:

-2lnL for Full Model 2344.609
 -2lnL for Intercept-Only Model 2398.103
 Chi-Square for Testing Intercept-Only Model 53.494
 Degrees of Freedom 3

PR NOSAY on EDUCAT LEFTRIGH AGE GENDER:

-2lnL for Full Model 2344.591
 -2lnL for Intercept-Only Model 2398.103
 Chi-Square for Testing Intercept-Only Model 53.512
 Degrees of Freedom 4

Recall that chi-square is a test of the hypothesis that none of the covariates has any effect. This hypothesis is rejected for any equation with EDUCAT included. Note that chi-square increases considerably when LEFTRIGH is added to EDUCAT and when AGE is added to EDUCAT and LEFTRIGH but not when GENDER is added. The chi-square difference $48.638 - 40.136 = 8.502$ with one degree of freedom is a test of the hypothesis that LEFTRIGH has no effect, given that EDUCAT is included. This hypothesis is rejected. Thus, LEFTRIGH should be included with EDUCAT. Similarly, the chi-square difference $53.494 - 48.638 = 4.856$ with one degree of freedom is a test of the hypothesis that AGE has no effect, given that EDUCAT and LEFTRIGH are included. This hypothesis is also rejected (at the 5% level). Thus, AGE should be included with EDUCAT and LEFTRIGH. But one cannot reject the hypothesis that GENDER has no effect, given that EDUCAT, LEFTRIGH, and AGE are included in the equation because $53.512 - 53.494 = 0.018$ is not significant.

Table 13: Estimated Category Probabilities (probit)

Covariates				Probabilities			
Gender	LeftRight	Education	Age	AS	A	D	DS
1	2	1	30	0.026	0.138	0.622	0.215
1	2	1	60	0.018	0.114	0.611	0.257
1	2	3	30	0.004	0.041	0.483	0.472
1	2	3	60	0.003	0.031	0.440	0.527
1	8	1	30	0.014	0.095	0.595	0.296
1	8	1	60	0.010	0.076	0.570	0.345
1	8	3	30	0.002	0.024	0.401	0.573
1	8	3	60	0.001	0.017	0.356	0.626
2	2	1	30	0.025	0.136	0.621	0.218
2	2	1	60	0.018	0.112	0.610	0.260
2	2	3	30	0.004	0.040	0.480	0.476
2	2	3	60	0.002	0.030	0.437	0.530
2	8	1	30	0.013	0.093	0.594	0.299
2	8	1	60	0.009	0.074	0.568	0.348
2	8	3	30	0.002	0.023	0.398	0.577
2	8	3	60	0.001	0.017	0.353	0.629

Using (22) and the estimated parameter values, one can compute estimated category probabilities for any specified set of covariate values. I illustrate this for the probit and logit regressions of NOSAY. Table 13 gives estimated probabilities for 16 different combinations of the four covariates for the probit model and Table 14 gives the same probabilities estimated under the logit model.

Table 13 shows that a young male with low education and "leftist" opinion is most likely to respond *Disagree* ($P = 0.622$) to the NOSAY statement. This may be contrasted with an old male with high education and "rightist" opinion whose most likely response is *Disagree Strongly* ($P = 0.626$). It is also seen that any person is more likely to respond *Disagree* or *Disagree Strongly* than *Agree* or *Agree Strongly* no matter what his characteristics are. The probability of an *Agree Strongly* response is very small for all types of persons. Table 14 shows very similar probabilities, but note that all probabilities for *Agree Strongly* are larger than the corresponding probabilities in Table 13 and most of the probabilities for *Disagree Strongly* are larger in Table 14 than in Table 13. This is in line with the remark made earlier that the logit model gives more probability to the outer categories than the probit model.

11. Probit and Logit Regression of All Efficacy Variables

To analyze all the six efficacy variables jointly with the four covariates, just replace the PR line in **ORD52.PRL** with (see **ORD53.PRL** where the Standard Parameterization is used):

```
PR NOSAY on GENDER - AGE
PR VOTING on GENDER - AGE
PR COMPLEX on GENDER - AGE
PR NOCARE on GENDER - AGE
PR TOUCH on GENDER - AGE
PR INTEREST on GENDER - AGE
```

A slight editing of the output file **ORD53.OUT** gives the following estimated probit regressions.

```
NOSAY = 0.0 + 0.00899*GENDER + 0.0424*LEFTRIGH + 0.360*EDUCAT + 0.00453*AGE +
Error, R2 = 0.059
Standerr      (0.0675)      (0.0184)      (0.0519)      (0.00207)
Z-values      0.133      2.308      6.946      2.187
P-values      0.894      0.021      0.000      0.029

VOTING = - 0.0344*GENDER - 0.0217*LEFTRIGH + 0.447*EDUCAT - 0.00634*AGE +
Error
              (0.0667) (0.0181)      (0.0515)      (0.00205)
              -0.516  -1.195      8.673      -3.096

COMPLEX = - 0.212*GENDER - 0.0233*LEFTRIGH + 0.494*EDUCAT + 0.000881*AGE +
Error
              (0.0678) (0.0184)      (0.0525)      (0.00207)
              -3.135  -1.26      9.402      0.425
```

Table 14: Estimated Category Probabilities (logit)

Covariates				Probabilities			
Gender	LeftRight	Education	Age	AS	A	D	DS
1	2	1	30	0.032	0.124	0.638	0.205
1	2	1	60	0.026	0.102	0.626	0.246
1	2	3	30	0.009	0.039	0.465	0.486
1	2	3	60	0.007	0.031	0.417	0.545
1	8	1	30	0.020	0.081	0.599	0.300
1	8	1	60	0.016	0.066	0.567	0.351
1	8	3	30	0.005	0.024	0.359	0.611
1	8	3	60	0.004	0.019	0.311	0.665
2	2	1	30	0.033	0.125	0.638	0.204
2	2	1	60	0.026	0.103	0.626	0.245
2	2	3	30	0.009	0.039	0.466	0.485
2	2	3	60	0.007	0.032	0.418	0.544
2	8	1	30	0.020	0.081	0.600	0.299
2	8	1	60	0.016	0.066	0.568	0.350
2	8	3	30	0.005	0.024	0.360	0.610
2	8	3	60	0.004	0.019	0.312	0.664

$\text{NOCARE} = -0.0402 \cdot \text{GENDER} + 0.0240 \cdot \text{LEFTRIGH} + 0.371 \cdot \text{EDUCAT} + 0.00288 \cdot \text{AGE} + \text{Error}$
 (0.0669) (0.0182) (0.0514) (0.00205)
 -0.600 1.319 7.219 1.407

$\text{TOUCH} = 0.0382 \cdot \text{GENDER} + 0.0118 \cdot \text{LEFTRIGH} + 0.290 \cdot \text{EDUCAT} + 0.00540 \cdot \text{AGE} + \text{Error}$
 (0.0676) (0.0184) (0.0516) (0.00207)
 0.565 0.643 5.632 2.604

$\text{INTEREST} = 0.0316 \cdot \text{GENDER} - 0.00604 \cdot \text{LEFTRIGH} + 0.249 \cdot \text{EDUCAT} + 0.00467 \cdot \text{AGE} + \text{Error}$
 (0.0672) (0.0183) (0.0511) (0.00206)
 0.470 -0.331 4.864 2.266

Thus,

- GENDER has a significant effect only for COMPLEX.
- LEFTRIGH is significant only for NOSAY.
- EDUCAT is significant for all the ordinal variables.
- AGE is significant for NOSAY, VOTING, TOUCH, and INTEREST. Note that the effect of AGE on VOTING is negative.

The output file also gives the following information about the fit of the probit regressions.

PR NOSAY on GENDER - AGE:

-2lnL for Full Model	2344.591
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	53.512
Degrees of Freedom	4

PR VOTING on GENDER - AGE:

-2lnL for Full Model	2470.648
-2lnL for Intercept-Only Model	2577.465
Chi-Square for Testing Intercept-Only Model	106.818
Degrees of Freedom	4

PR COMPLEX on GENDER - AGE:

-2lnL for Full Model	2284.676
-2lnL for Intercept-Only Model	2393.401
Chi-Square for Testing Intercept-Only Model	108.725
Degrees of Freedom	4

PR NOCARE on GENDER - AGE:

-2lnL for Full Model	2401.907
-2lnL for Intercept-Only Model	2455.675
Chi-Square for Testing Intercept-Only Model	53.767
Degrees of Freedom	4

PR TOUCH on GENDER - AGE:

-2lnL for Full Model	2249.115
-2lnL for Intercept-Only Model	2284.420
Chi-Square for Testing Intercept-Only Model	35.305
Degrees of Freedom	4

PR INTEREST on GENDER - AGE:

-2lnL for Full Model	2338.651
-2lnL for Intercept-Only Model	2364.814
Chi-Square for Testing Intercept-Only Model	26.163
Degrees of Freedom	4

The second line for each regression gives a deviance but since we have no base model to compare it with this does not provide any information about whether the probit model fits the data or not. The third and fourth lines give a chi-square test of the hypothesis that all regression coefficients are zero. It is seen that this hypothesis is rejected for all ordinal variables. This is as it should be.

For comparison, I give the fit statistics for the logit regressions obtained by putting LR instead of PR in **ORD53.PRL**, see file **ORD53A.PRL**.

LR NOSAY on GENDER - AGE:

-2lnL for Full Model	2340.790
-2lnL for Intercept-Only Model	2398.103
Chi-Square for Testing Intercept-Only Model	57.313
Degrees of Freedom	4

LR VOTING on GENDER - AGE:

-2lnL for Full Model	2467.314
-2lnL for Intercept-Only Model	2577.465
Chi-Square for Testing Intercept-Only Model	110.152
Degrees of Freedom	4

LR COMPLEX on GENDER - AGE:

-2lnL for Full Model	2279.764
-2lnL for Intercept-Only Model	2393.401
Chi-Square for Testing Intercept-Only Model	113.637
Degrees of Freedom	4

LR NOCARE on GENDER - AGE:

-2lnL for Full Model	2392.594
-2lnL for Intercept-Only Model	2455.675
Chi-Square for Testing Intercept-Only Model	63.081
Degrees of Freedom	4

LR TOUCH on GENDER - AGE:

-2lnL for Full Model	2244.766
-2lnL for Intercept-Only Model	2284.420
Chi-Square for Testing Intercept-Only Model	39.653
Degrees of Freedom	4

LR INTEREST on GENDER - AGE:

-2lnL for Full Model	2334.276
-2lnL for Intercept-Only Model	2364.814
Chi-Square for Testing Intercept-Only Model	30.538
Degrees of Freedom	4

Does the logit model fit better than the probit model? The answer is Yes, Yes, Yes, Yes, Yes, and Yes. The two models have the same number of parameters but the deviance is smaller for the logit model than for the probit model for all variables. Take NOSAY, for example. The difference in deviance is $2344.591 - 2340.790 = 3.801$. Note that one can obtain the same number as the difference between the two chi-squares in the reverse order: $57.313 - 53.512 = 3.801$.

12. Estimating the Joint Covariance Matrix

To estimate the joint covariance matrix of the continuous variables underlying the ordinal variables and the covariates as defined in Section 6, one must use a `Fixedvariables` command (or `FI` command for short), see Jöreskog & Sörbom (1999a, pp. 180-183). Instead of `Fixedvariables` one can write `Covariates`. In addition to all probit regressions, these commands give estimates of the conditional covariance matrix and the joint unconditional covariance matrix as defined in Section 6. File **ORD54.PRL** illustrates this using the Standard Parameterization. It also shows how one can obtain the asymptotic covariance matrix of the joint unconditional covariance matrix. File **ORD54.PRL** is

```
Computing Covariance Matrix
SY=USA.LSF
Fixedvariables: GENDER - AGE
Output MA=CM CM=USA.CM AC=USA.ACC WP
```

All variables specified on the `Covariates:` line are automatically treated as continuous variables. All other variables are assumed to be ordinal.

The output file **ORD54.OUT** gives the conditional covariance matrix as

Conditional Covariance Matrix

	NOSAY	VOTING	COMPLEX	NOCARE	TOUCH	INTEREST
	-----	-----	-----	-----	-----	-----
NOSAY	1.000					
VOTING	0.284 (0.034) 8.312	1.000				
COMPLEX	0.270 (0.035) 7.746	0.204 (0.035) 5.766	1.000			
NOCARE	0.567 (0.027) 21.089	0.223 (0.035) 6.412	0.379 (0.032) 11.720	1.000		
TOUCH	0.367 (0.033) 11.106	0.206 (0.035) 5.824	0.274 (0.035) 7.874	0.637 (0.024) 26.169	1.000	
INTEREST	0.460 (0.030) 15.117	0.200 (0.035) 5.677	0.305 (0.034) 8.963	0.657 (0.023) 28.183	0.674 (0.023) 29.221	1.000

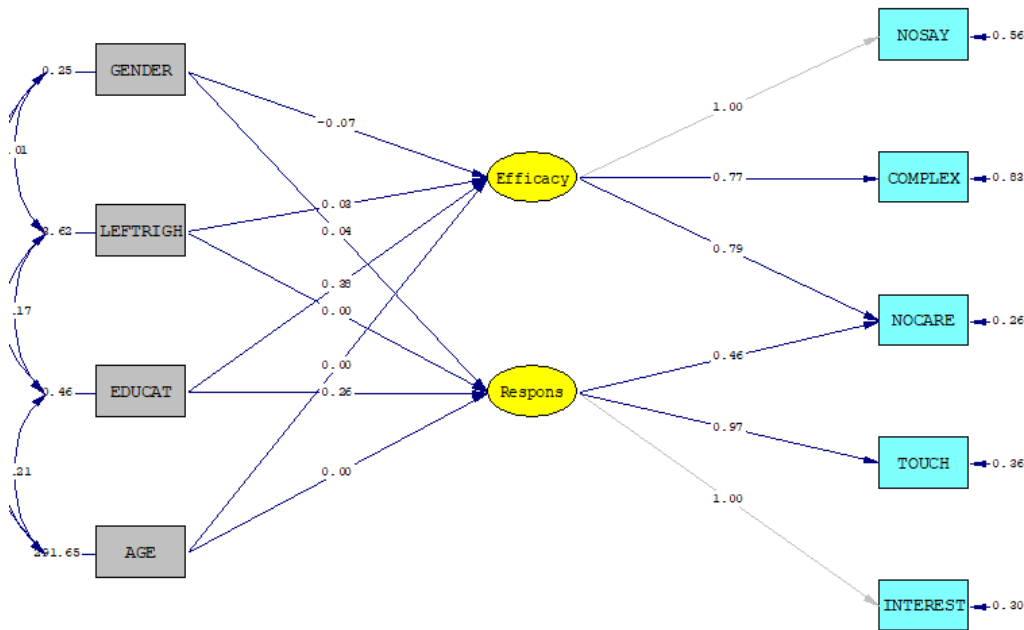
In this case, when the Standard Parameterization is used, this is the correlation matrix of the error terms. All correlations are highly significant. This means that the covariates alone do not account for the correlations of the ordinal variables (or more correctly the variables underlying the ordinal variables). This is not surprising since we know from the second example in this set that we need the latent variables *Efficacy* and *Respons* to account for these correlations. In Section 13 I will use these latent variables as well.

The output **ORD54.OUT** also gives the joint covariance matrix of the variables underlying the ordinal variables and the covariates. This is too large to list here. It is saved in the file **USA.CM** and its asymptotic covariance matrix is saved in the file **USA.ACC**. The covariance matrix **USA.CM** is an unconstrained covariance matrix just as a sample covariance matrix for continuous variables. It can therefore be used for modeling in LISREL just as if all variables were continuous. The only restriction is that the covariates must not be treated as indicators of latent variables. In LISREL, one can estimate the model either by WLS using the inverse of **USA.ACC** as a weight matrix or by ML using **USA.ACC** to correct standard errors and chi-square for non-normality.

In the second example, I used the ordinal *Efficacy* variables to establish a measurement model for the two latent variables *Efficacy* and *Respons*. Now I will investigate to what extent the covariates affect these two latent variables. To investigate this, one can use a MIMIC model described in Section 13.

13. A MIMIC Model for Efficacy and Respons

The idea of a MIMIC model is that a set of possibly explanatory variables (covariates) affects latent variables which are indicated by other observed variables, in this case ordinal variables. Thus there are multiple indicators and multiple causes of latent variables, see Jöreskog & Goldberger (1975). For examples of MIMIC models with continuous indicators see Jöreskog & Sörbom (1999b). The MIMIC model considered here is shown in Fig. 10.



Chi-Square=94.02, df=15, P-value=0.00000, RMSEA=0.070

Figure 10: MIMIC Model for Efficacy and Respons

A SIMPLIS command file for estimating the model in Fig. 10 is **ORD55.SPL**:

```
MIMIC Model
Observed Variables: NOSAY VOTING COMPLEX NOCARE TOUCH INTEREST
GENDER LEFTRIGH EDUCAT AGE
Covariance Matrix from File USA.CM
Asymptotic Covariance Matrix from File USA.ACC
Sample Size: 1076
Latent Variables: Efficacy Respons
Relationships:
  NOSAY COMPLEX NOCARE = Efficacy
  NOCARE TOUCH INTEREST = Respons
  NOSAY = 1*Efficacy
  INTEREST = 1*Respons
  Efficacy Respons = GENDER LEFTRIGH EDUCAT AGE
Let the errors of Efficacy and Respons correlate
Path Diagram
End of Problem
```


The output gives the structural equations as

```

Efficacy = - 0.0717*GENDER + 0.0251*LEFTRIGH + 0.381*EDUCAT + 0.00225*AGE, Errorvar.=
0.435 , R² = 0.131
Standerr      (0.0497)          (0.0135)          (0.0406)          (0.00151)
(0.0411)
Z-values      -1.441            1.860            9.392            1.485
10.572
P-values      0.149            0.063            0.000            0.137            0.000

Respons = 0.0352*GENDER + 0.00277*LEFTRIGH + 0.257*EDUCAT + 0.00476*AGE, Errorvar.=
0.701 , R² = 0.0430
Standerr (0.0567)      (0.0154)      (0.0431)      (0.00174)      (0.0448)
Z-values  0.621        0.180        5.962        2.738        15.633
P-values  0.534        0.858        0.000        0.006        0.000

```

which shows that GENDER has a significant effect on *Efficacy* and EDUCAT has significant effects on both *Efficacy* and *Respons*. LEFTRIGH and AGE have no significant effects on either of the latent variables. The fact that they are non-significant does not mean they do not exist, only that the sample size is not sufficiently large to make them significant.

The model fits the data reasonably well as judged by the following fit statistics. For this conclusion I use the information about RMSEA and the guidelines of Browne & Cudeck (1993).

Goodness-of-Fit Statistics

Degrees of Freedom for (C1)-(C2)	15
Maximum Likelihood Ratio Chi-Square (C1)	94.022 (P = 0.0000)
Browne's (1984) ADF Chi-Square (C2_NT)	92.916 (P = 0.0000)
Estimated Non-centrality Parameter (NCP)	79.022
90 Percent Confidence Interval for NCP	(52.089 ; 113.457)
Minimum Fit Function Value	0.0875
Population Discrepancy Function Value (F0)	0.0735
90 Percent Confidence Interval for F0	(0.0485 ; 0.106)
Root Mean Square Error of Approximation (RMSEA)	0.0700
90 Percent Confidence Interval for RMSEA	(0.0568 ; 0.0839)
P-Value for Test of Close Fit (RMSEA < 0.05)	0.00692