

Structural equation models for regression with measurement error: Just enough theory (Draft Two)

Jerry Brunner and Christine Lim
Department of Statistics, University of Toronto
100 St. George St.
Toronto M5G 3G3, Canada

November 16, 2007

Abstract

This article shows how structural equation modelling methods may be used to carry out a valid regression analysis when independent variables are measured with error. An essential part of the process is to overcome the problem of model identification. A general solution, called the “double measurement design,” is described; this involves measuring each variable twice, in such a way that errors of measurement are independent between occasions. When data are collected according to the double measurement design, model identification is guaranteed, and the data analyst need not struggle with mathematical details.

Keywords: Errors in variables, Model identification, Tau-equivalent measures, Double measurement.

Introduction

In a survey, suppose that a respondent’s annual income is “measured” by simply asking how much he or she earned last year. Will this measurement be completely accurate? Of course not. Some people will lie, some will forget and give a reasonable guess, and still others will suffer from legitimate confusion about what constitutes income. Even physical variables like height, weight and blood pressure are subject to some inexactness of measurement, no matter how skilled the personnel doing the measuring. In fact, very few of the variables in the typical data set are measured completely without error. An exception might be something like the amount of drug administered in a clinical trial. Here, laboratory procedures guarantee that for all practical purposes, the amount of drug a subject receives is exactly what you think it is. But in general, if a variable is simply measured rather than being experimentally manipulated, there is usually at least a little bit of measurement error.

Random variables that cannot be directly observed are called *latent variables*. The ones we can observe are sometimes called “manifest,” but in this paper they will be called “observed” or “observable,” which is also a common usage. Upon reflection, it is clear that most of the time, we are interested in relationships among latent variables, but at best our data consist only of their imperfect, observable counterparts. One is reminded of the allegory of the cave in Plato’s *Republic*, where human beings are compared to prisoners in a cave, with their heads chained so that they can only look at a wall. Behind them is a fire, which casts flickering shadows on the wall. They cannot observe reality directly; all they can see are the shadows.

In ordinary least-squares regression, the only latent variable is the error term. Measurement error in the dependent variable can perhaps be absorbed into the error term, but there is no provision for measurement error in the independent variables. Unfortunately, when independent variables are measured with error, the results can be disastrous. Estimated regression coefficients are biased even as the sample size approaches infinity, and Type I error rates can be seriously inflated.

This has been known for a long time. The alarm about biased regression coefficients was sounded by Stouffer (1936), and by the seventh edition of *Statistical methods for research workers*, Fisher (1938) was warning scientists about the problem. For a modern and readable discussion of what happens to ordinary least-squares regression when measurement error is ignored, the classic article by Cochran (1968) is an excellent source. Fuller (1987) provides an authoritative treatment of regression models that incorporate measurement error; also see Cheng and Van Ness (1999). And the classical structural equation models (for example Goldberger and Duncan, 1973; Jöreskog, 1978; McArdale, 1980; McDonald, 1978; Bentler and Weeks, 1980; Bollen, 1989) include regression with and without measurement error as special cases.

Nevertheless, few regression texts outside Econometrics provide guidance about what to do when the independent variables are measured with error. The present article attempts to fill this gap. It uses language and notation associated with the LISREL structural equation model (Jöreskog, 1978; Bollen, 1989) rather than the arguably more sophisticated approach of Fuller (1987), in order to be accessible to advanced undergraduates in Statistics. Another advantage of the structural equation modelling approach is that high-quality commercial software is available. SAS `proc calis` (SAS Institute, 1999) is available in many academic environments, and LISREL (Jöreskog and Sörbom, 1996) and AMOS (Arbuckle, 2006) are excellent programs with free student versions. There is also a structural equation modelling package for R (Fox, 2006).

Here is the plan of the paper. Section 1 presents almost the simplest possible regression model with measurement error. There is one independent variable, no intercept, additive measurement error, and everything is normally distributed. We will see that even in this case, the model parameters cannot be successfully estimated from the data. The problem is *model identification*. When a statistical model is not identified, it is impossible to recover the parameters even from an infinite amount of data.

Section 2 discusses model identification, and arrives at a well-known principle that

applies to all structural equation models, including models of regression with measurement error. The principle is this. The mean and covariance matrix of the observable variables are always functions of the model parameters. If the model parameters are also functions of the mean and covariance matrix, then those parameters are identified.

Section 3 describes a general solution of the identification problem for regression with measurement error: the *double measurement design*. This consists of measuring all the variables twice, preferably on two different occasions, with different measurement procedures. If this can be done in such a way that the errors of measurement on the two occasions are independent, then model identification is taken care of automatically, and the analysis can proceed in a routine manner.

The double measurement design is similar to the idea of “tau-equivalent measures” (for example Bollen, 1987, p. 208, or cite Lord and Novick?), except that all measurement errors need not be independent. In fact, a very desirable feature of the double measurement design is that while errors of measurement from different measurement procedures must be independent, errors of measurement from the same measurement procedure are allowed to be correlated. For example, one should always expect correlated measurement errors for self-report data; these would arise from consistent individual differences in style of responding to questionnaires and in desire to make a favorable impression. And when measurement errors are correlated, adopting a model where they are uncorrelated can have effects that are just as bad as ignoring measurement error altogether.

The double measurement model of Section 3 employs the classical Structural Equation Modelling trick of “centering” all the variables by subtracting off the means, and then conducting the analysis under the assumption that all expected values are zero. In Section 4, the model is expanded to include intercepts. But in most cases this just makes the model parameters harder to identify, and does not providing any additional information about the relationship between the independent and dependent variables. The final conclusion is that most of the time, including intercepts is not worth the extra trouble.

In the development of this theory, assuming multivariate normality simplifies the exposition but is not really necessary. In Section 5, the normal assumption is relaxed. For independent variables that are measured without error (for example, the dummy variables for factors that are experimentally manipulated), the distribution does not matter at all. For independent variables that are measured with error, the double measurement design guarantees identification of a necessary *function* of the parameters of a distribution-free model.

The double measurement design also points to Method of Moments estimators of the regression coefficient that are consistent and asymptotically normal, by a straightforward application of the Central Limit Theorem. This would provide the basis for a full set of large-sample tests and confidence intervals, but it is unnecessary to go there. In fact, the estimators and tests based on a multivariate normal assumption enjoy robustness properties that make them superior to one method (the weighted least-squares approach of Browne, 1984) that was specifically designed to avoid the assumption of normality. They are probably also superior to the methods suggested by the double measurement

design, which are very similar to Browne's.

The moral of the story is comforting in its simplicity. For data that are collected according to the double measurement recipe, just fit a classical structural model with no intercepts and everything normally distributed; this is close to the default settings of most available software.

1 Regression through the origin with one independent variable

Even in the simplest case, when we try to incorporate measurement error into a regression model, we immediately encounter a technical difficulty: model identification. In a simple regression, suppose the dependent variable is related to an independent variable. We can observe the dependent variable, but not the actual value of the independent variable. All we can see is the independent variable plus a piece of random noise.

Independently for $i = 1, \dots, n$, let

$$\begin{aligned} Y_i &= \gamma \xi_i + \zeta_i \\ X_i &= \xi_i + \delta_i, \end{aligned} \tag{1}$$

where ξ_i , ζ_i and δ_i are independent normal random variables with expected value zero, $Var(\xi_i) = \phi$, $Var(\zeta_i) = \psi$, and $Var(\delta_i) = \theta_\delta$. The regression coefficient γ is a fixed constant. The notation here is taken from the LISREL structural equation model (Jöreskog, 1978; Bollen, 1989) for compatibility with later parts of this paper, and because familiarity with this notation will make it easier for students to use structural equation modelling software.

Data from Model (1) are just the pairs (X_i, Y_i) for $i = 1, \dots, n$. The true independent variable ξ_i is a latent variable whose value cannot be known exactly. The model implies that the (X_i, Y_i) are independent bivariate normal with mean zero and covariance matrix

$$\Sigma = \begin{bmatrix} \phi + \theta_\delta & \gamma\phi \\ \gamma\phi & \gamma^2\phi + \psi \end{bmatrix}. \tag{2}$$

A multivariate normal distribution with mean zero is completely characterized by its covariance matrix, so even an infinite amount of data can only tell us the three unique values in the matrix Σ . But there are four parameters in the model: γ , ϕ , ψ and θ_δ . Recovering all four parameters from the unique elements of Σ amounts to solving three equations in four unknowns — an impossibility. Maximum likelihood estimation will fail, with a non-unique maximum at an infinite number of points along a curve in four dimensions.

The problem is that Model (1) is not uniquely identified in the model parameters. The concept of model identification is unfamiliar to most students, because typically (except in the case of exploratory factor analysis) we present them with statistical models that

are nicely identified, and the issue does not arise. Thus, a general discussion of model identification may be helpful.

2 Model identification

Suppose we have a vector of observable data $\mathbf{D} = (D_1, \dots, D_n)$, and a statistical model (a set of assertions implying a probability distribution) for \mathbf{D} . The model depends on a parameter θ , which is usually a vector. If the probability distribution of \mathbf{D} corresponds uniquely to θ , then we say that the model is *identified*. But if any two different parameter values yield the same probability distribution, then the model is not identified. In this case, the data cannot be used to decide between the two parameter values, and standard methods of parameter estimation will fail. Even an infinite amount of data cannot tell you the true parameter values.

In Model (1), $\theta = (\gamma, \phi, \psi, \theta_\delta)$, $D_i = (X_i, Y_i)$, and the probability distribution of \mathbf{D} is completely determined by Σ . The two variances and one covariance in Σ cannot correspond uniquely to the four elements of θ , so the model is not identified. To really nail it down, the two distinct parameter values $\theta_1 = (1, 2, 4, 1)$ and $\theta_2 = (2, 1, 1, 2)$ both yield

$$\Sigma = \begin{bmatrix} 3 & 2 \\ 2 & 6 \end{bmatrix},$$

and thus the same distribution of the sample data. The clearest way to prove a model is non-identified is with a simple numerical example like this, but frequently other arguments are more convenient.

It is sometimes useful to distinguish among three kinds of model identification. In *pointwise* identification, a model is said to be identified at a point if no other point in the parameter space yields the same probability distribution. If the model is identified at every point in the parameter space, it is said to be *globally* identified. The model is called *locally* identified at a point θ_0 if that point is surrounded by a neighborhood where none of the other points in the neighborhood yield the same probability distribution as θ_0 . A globally identified model is locally identified at every point, but the converse is not true. In this paper, “identification” refers to global identification.

It is possible for certain *functions* of the parameter vector to be identified, even when the entire model is not. If full knowledge of the probability distribution of \mathbf{D} implies knowledge of some function of θ , then that function is said to be identified, and consistent estimation of it is a possibility. For example, let D_1, \dots, D_n be i.i.d. Poisson random variables with mean $\lambda_1 + \lambda_2$, where $\lambda_1 > 0$ and $\lambda_2 > 0$. The parameter is the pair $\theta = (\lambda_1, \lambda_2)$. The model is not identified because any pair of λ values satisfying $\lambda_1 + \lambda_2 = c$ will produce exactly the same probability distribution. Notice also how maximum likelihood estimation will fail in this case; the likelihood function will have a ridge, a non-unique maximum along the line $\lambda_1 + \lambda_2 = \bar{D}$, where \bar{D} is the sample mean. The function $\lambda = \lambda_1 + \lambda_2$, of course, is identified.

The Normal distribution Suppose we have a random sample D_1, \dots, D_n from a normal distribution with parameters μ and σ^2 . This model must be identified because we use it all the time, but how can one prove it? An indirect way is to observe that consistent estimators of μ and σ^2 are available, and this would not be possible if the model were not identified. But it is more instructive to use the definition. To clarify the ideas, we shall be very explicit.

The objective to show that the parameter $\theta = (\mu, \sigma^2)$ is a function of the distribution of the data vector $\mathbf{D} = (D_1, \dots, D_n)$. By “distribution,” we mean the joint cumulative distribution function, which in this case may be written

$$F_{\mathbf{D}}(\mathbf{d}; \theta) = \prod_{i=1}^n F(d_i; \theta) = \prod_{i=1}^n \int_{-\infty}^{d_i} \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x_i - \mu)^2\right] dx_i.$$

For independent and identically distributed data, we only need the marginal distribution of one observation in order to check model identification. The marginal distribution of D_1 will do. It is a function of the joint distribution.

$$g_1(F_{\mathbf{D}}) = \lim_{d_2 \rightarrow \infty} \dots \lim_{d_n \rightarrow \infty} F_{\mathbf{D}}(\mathbf{d}; \theta) = F(d_1; \theta) \prod_{i=2}^n \lim_{d_i \rightarrow \infty} F(d_i; \theta) = F(d_1; \theta)$$

Now we take a function g_2 of the marginal distribution. This function maps a space of cumulative distribution functions into \mathfrak{R}^2 .

$$g_2(F) = g_2(g_1(F_{\mathbf{D}})) = \begin{bmatrix} \int_{-\infty}^{\infty} x F'(x; \theta) dx \\ \int_{-\infty}^{\infty} x^2 F'(x; \theta) dx \end{bmatrix} = \begin{bmatrix} E(D_1) \\ E(D_1^2) \end{bmatrix}$$

And finally, one last function $g_3 : \mathfrak{R}^2 \rightarrow \mathfrak{R}^2$:

$$g_3(g_2(g_1(F_{\mathbf{D}}))) = \begin{bmatrix} E(D_1) \\ E(D_1^2) - [E(D_1)]^2 \end{bmatrix} = \begin{bmatrix} \mu \\ \sigma^2 \end{bmatrix}$$

This last composite function expresses the parameter as a function of the distribution of the data, proving model identification. It is usually unnecessary to go into such detail. In general, if the parameters of a statistical model can be recovered from the moments of the model distribution, then the model is identified. If a *function* of the parameters can be recovered from the moments, then that function of the parameters is identified. Identification of most common probability models (including the multivariate normal) follows in this way.

Back to regression with measurement error Classical structural equation models, including models for regression with measurement error, are based on systems of simultaneous linear equations. Assuming simple random sampling from a large population, the observable data are independent and identically distributed, with a mean vector $\boldsymbol{\mu}$ and a covariance matrix $\boldsymbol{\Sigma}$ that may be written as functions of the model parameters in a

straightforward way. If it is possible to solve uniquely for a given model parameter in terms of the elements of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, then that parameter is a function of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, which in turn are functions of the probability distribution of the data. A function of a function is a function, and so the parameter is a function of the probability distribution of the data. Hence, it is identified.

To summarize, we have arrived at the standard way to check model identification for any linear simultaneous equation model, not just measurement error regression. *First, calculate the expected value and covariance matrix of the observable data, as a function of the model parameters. If it is possible to solve uniquely for the model parameters in terms of the means, variances and covariances of the observable data, then the model parameters are identified.* If all the random vectors in the model are multivariate normal, this condition is necessary as well as sufficient.

Example: Instrumental variables In a model like (1), suppose that we have access to data for another two variables that depend on the latent independent variable ξ . Our main interest is still in Y ; the other two are called *instrumental* variables because they are just tools for obtaining an identified model.

Here is the expanded version of Model (1). The original dependent variable Y is now called Y_1 . Following the usual convention in structural equation modelling, the subscript i has been omitted to reduce notational clutter. The model is presented for a single observation, and implicitly everything is independent and identically distributed, for $i = 1, \dots, n$.

$$\begin{aligned} X &= \xi + \delta \\ Y_1 &= \gamma_1 \xi + \zeta_1 \\ Y_2 &= \gamma_2 \xi + \zeta_2 \\ Y_3 &= \gamma_3 \xi + \zeta_3, \end{aligned} \tag{3}$$

where δ , ξ , ζ_1 , ζ_2 and ζ_3 are all independent, $Var(\xi) = \phi$, $Var(\zeta_1) = \psi_1$, $Var(\zeta_2) = \psi_2$, $Var(\zeta_3) = \psi_3$, $Var(\delta) = \theta_\delta$, all expected values are zero, and the regression coefficients γ_1 , γ_2 and γ_3 are fixed constants.

Writing the vector of observable data (for subject i) as $\mathbf{D} = (X, Y_1, Y_2, Y_3)'$, elements of the covariance matrix $\boldsymbol{\Sigma}$ may be obtained by elementary one-variable calculations, like $Var(X) = Var(\xi + \delta) = Var(\xi) + Var(\delta) = \phi + \theta_\delta$, and

$$\begin{aligned} Cov(X, Y_1) &= E(X, Y_1) = E([\xi + \delta][\gamma_1 \xi + \zeta_1]) = E(\gamma_1 \xi^2 + \xi \zeta_1 + \gamma_1 \delta \xi + \delta \zeta_1) \\ &= \gamma_1 E(\xi^2) + E(\xi \zeta_1) + \gamma_1 E(\delta \xi) + E(\delta \zeta_1) \\ &= \gamma_1 Var(\xi) + E(\xi)E(\zeta_1) + \gamma_1 E(\delta)E(\xi) + E(\delta)E(\zeta_1) \\ &= \gamma_1 \phi \end{aligned}$$

In this way, we obtain

$$\Sigma = \begin{bmatrix} \phi + \theta_\delta & \gamma_1\phi & \gamma_2\phi & \gamma_3\phi \\ \gamma_1^2\phi + \psi_1 & \gamma_1\gamma_2\phi & \gamma_1\gamma_3\phi & \\ \gamma_2^2\phi + \psi_2 & \gamma_2\gamma_3\phi & & \\ \gamma_3^2\phi + \psi_3 & & & \end{bmatrix}. \quad (4)$$

To prove model identification, we need to solve for the model parameters in terms of Σ . Denote the i, j element of Σ by σ_{ij} . The task is to solve the following ten equations in eight unknowns

$$\begin{aligned} \sigma_{11} &= \phi + \theta_\delta \\ \sigma_{12} &= \gamma_1\phi \\ \sigma_{13} &= \gamma_2\phi \\ \sigma_{14} &= \gamma_3\phi \\ \sigma_{22} &= \gamma_1^2\phi + \psi_1 \\ \sigma_{23} &= \gamma_1\gamma_2\phi \\ \sigma_{24} &= \gamma_1\gamma_3\phi \\ \sigma_{33} &= \gamma_2^2\phi + \psi_2 \\ \sigma_{34} &= \gamma_2\gamma_3\phi \\ \sigma_{44} &= \gamma_3^2\phi + \psi_3 \end{aligned} \quad (5)$$

for $\phi, \theta_\delta, \gamma_1, \gamma_2, \gamma_3, \psi_1, \psi_2,$ and ψ_3 .

The fact that there are more equations than unknowns does not guarantee the existence of a unique solution; it merely tells us that a unique solution is possible. Suppose that γ_2 and γ_3 are both non-zero. This is reasonable, because to be useful, the instrumental dependent variables must have some relationship to the independent variable. In this case,

$$\frac{\sigma_{13}\sigma_{14}}{\sigma_{34}} = \frac{\gamma_2\gamma_3\phi^2}{\gamma_2\gamma_3\phi} = \phi. \quad (6)$$

Then, simple substitutions allow us to solve for the rest of the parameters, yielding the complete solution

$$\begin{aligned} \phi &= \frac{\sigma_{13}\sigma_{14}}{\sigma_{34}} \\ \theta_\delta &= \sigma_{11} - \frac{\sigma_{13}\sigma_{14}}{\sigma_{34}} \\ \gamma_1 &= \frac{\sigma_{12}\sigma_{34}}{\sigma_{13}\sigma_{14}} \end{aligned} \quad (7)$$

$$\begin{aligned}
\gamma_2 &= \frac{\sigma_{34}}{\sigma_{14}} \\
\gamma_3 &= \frac{\sigma_{34}}{\sigma_{13}} \\
\psi_1 &= \sigma_{22} - \frac{\sigma_{12}^2 \sigma_{34}}{\sigma_{13} \sigma_{14}} \\
\psi_2 &= \sigma_{33} - \frac{\sigma_{13} \sigma_{34}}{\sigma_{14}} \\
\psi_3 &= \sigma_{44} - \frac{\sigma_{14} \sigma_{34}}{\sigma_{13}}
\end{aligned}$$

This proves model identification. The solution is thorough but somewhat tedious, even for this simple example. The student may wonder how much work really needs to be shown. I would suggest showing the calculations leading to the covariance matrix (4), saying “Denote the i, j element of Σ by σ_{ij} ,” skipping the system of equations (5) because they are present in (4), and showing the solution for ϕ in (6), *including* the stipulation that γ_2 and γ_3 are both non-zero. Then, instead of the explicit solution (7), write something like

$$\begin{aligned}
\theta_\delta &= \sigma_{11} - \phi & (8) \\
\gamma_1 &= \frac{\sigma_{12}}{\phi} \\
\gamma_2 &= \frac{\sigma_{13}}{\phi} \\
\gamma_3 &= \frac{\sigma_{14}}{\phi} \\
\psi_1 &= \sigma_{22} - \gamma_1^2 \phi \\
\psi_2 &= \sigma_{33} - \gamma_2^2 \phi \\
\psi_3 &= \sigma_{44} - \gamma_3^2 \phi
\end{aligned}$$

Notice how once we have solved for a model parameter, we use it to solve for other parameters without explicitly substituting in terms of σ_{ij} . The objective is to prove that a unique solution exists by showing how to get it. An exact statement of the solution is not necessary.

Two additional comments are in order. First, this model had no intercepts, and the random variables all had expected value zero. This is typical of the classical structural equation models, in which inference is based solely on the sample covariance matrix and not the means. One speaks of “centering” all the variables by subtracting off the sample means (for example Bollen, 1989). For large samples, this is almost the same as subtracting off the population means. Since all the confidence intervals and tests are based on large-sample theory anyway, no harm is done. Later, we shall consider models with intercepts.

A second comment is that even for the most complex models, proving model identification as in the preceding example involves only elementary mathematics. But it can be long and messy, especially for models with lots of independent variables — and almost all real-life regressions have lots of independent variables. Furthermore, for a given data set, it is not always possible to come up with a realistic model that is identified. A sensible alternative is to plan the statistical analysis in advance, and to ensure model identification by collecting the right kind of data. The next section describes a way to do this. The key is to measure the independent variables twice, preferably using different methods or measuring instruments.

3 The double measurement design

For regression with measurement error, the model identification problem is solved if we measure all the variables on more than one occasion, in such a way that errors of measurement on different occasions are independent. We begin with a classical structural equation model in which all random variables have expected value zero and there no intercepts. In Section 4, the model is extended to include intercepts and non-zero expected values, but ultimately Model (9) below is recommended for most purposes.

For each of n independent observations, assume the following simultaneous equation model. Implicitly, all the random quantities involved have a subscript i , $i = 1, \dots, n$.

$$\begin{aligned}
 \boldsymbol{\eta} &= \boldsymbol{\Gamma}\boldsymbol{\xi} + \boldsymbol{\zeta} \\
 \mathbf{X}_1 &= \boldsymbol{\xi} + \boldsymbol{\delta}_1 \\
 \mathbf{X}_2 &= \boldsymbol{\xi} + \boldsymbol{\delta}_2, \\
 \mathbf{Y}_1 &= \boldsymbol{\eta} + \boldsymbol{\epsilon}_1 \\
 \mathbf{Y}_2 &= \boldsymbol{\eta} + \boldsymbol{\epsilon}_2,
 \end{aligned} \tag{9}$$

where

$\boldsymbol{\eta}$ is an $m \times 1$ random vector of latent dependent variables. Because m can be greater than one, the regression is multivariate.

$\boldsymbol{\Gamma}$ is an $m \times p$ matrix of unknown constants. These are the regression coefficients, with one row for each dependent variable and one column for each independent variable.

$\boldsymbol{\xi}$ is a $p \times 1$ random vector of latent independent variables, with expected value zero and variance-covariance matrix $\boldsymbol{\Phi}$, a $p \times p$ symmetric and positive definite matrix of unknown constants.

ζ is the error term of the latent regression. It is an $m \times 1$ random vector with expected value zero and variance-covariance matrix Ψ , an $m \times m$ symmetric and positive definite matrix of unknown constants.

\mathbf{X}_1 and \mathbf{X}_2 are $p \times 1$ observable random vectors, each representing ξ plus a different piece of random error.

δ_1 is the measurement error in \mathbf{X}_1 . It is a $p \times 1$ random vector of error terms, with expected value zero and variance-covariance matrix Θ_{δ_1} , a $p \times p$ symmetric and positive definite matrix of unknown constants.

δ_2 is the measurement error in \mathbf{X}_2 . It is a $p \times 1$ random vector of error terms, with expected value zero and variance-covariance matrix Θ_{δ_2} , a $p \times p$ symmetric and positive definite matrix of unknown constants.

\mathbf{Y}_1 and \mathbf{Y}_2 are $m \times 1$ observable random vectors, each representing η plus a different piece of random error.

ϵ_1 is the measurement error in \mathbf{Y}_1 . It is an $m \times 1$ random vector of error terms, with expected value zero and variance-covariance matrix Θ_{ϵ_1} , an $m \times m$ symmetric and positive definite matrix of unknown constants.

ϵ_2 is the measurement error in \mathbf{Y}_2 . It is a $p \times 1$ random vector of error terms, with expected value zero and variance-covariance matrix Θ_{ϵ_2} , an $m \times m$ symmetric and positive definite matrix of unknown constants.

ξ and ζ are independent of one another.

ξ and ζ are independent of (δ_1, ϵ_1) .

ξ and ζ are independent of (δ_2, ϵ_2) .

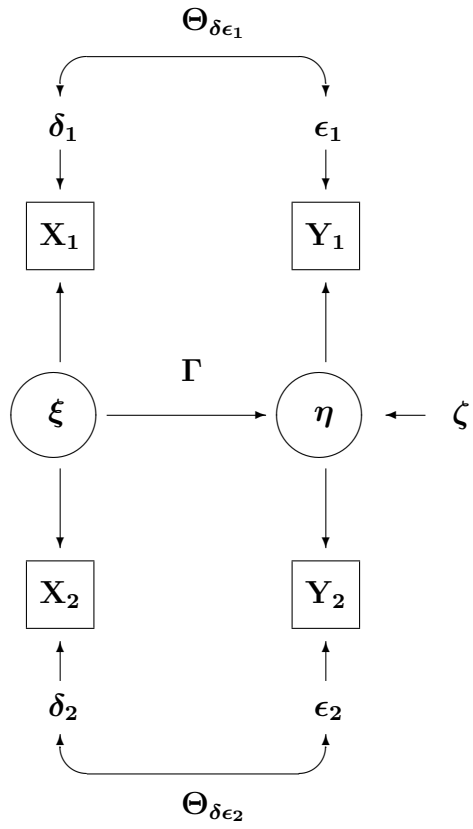
(δ_1, ϵ_1) are independent of (δ_2, ϵ_2) . However,

δ_1 and ϵ_1 need not be independent. $Cov(\delta_1, \epsilon_1) = \Theta_{\delta\epsilon_1}$, a $p \times m$ matrix of unknown constants.

δ_2 and ϵ_2 need not be independent. $Cov(\delta_2, \epsilon_2) = \Theta_{\delta\epsilon_2}$, a $p \times m$ matrix of unknown constants.

The model is depicted in Figure 1. It follows the usual conventions for path diagrams of structural equation models. Straight arrows go from *exogenous* variables (that is, independent variables, those on the right-hand side of equations) to *endogenous* variables (dependent variables, those on the left side). Correlations among exogenous variables are represented by two-headed curved arrows. Observable variables are enclosed by rectangles

Figure 1: The Double Measurement Model



or squares, while latent variables are enclosed by ellipses or circles. Error terms are not enclosed by anything.

The main idea of the Double Measurement Design is that every variable is measured by two different methods. Errors of measurement may be correlated within measurement methods, but not between methods. So for example, farmers who overestimate their number of pigs may also overestimate their number of cows. On the other hand, if the number of pigs is counted once by the farm manager at feeding time and on another occasion by a research assistant from an areal photograph, then it would be fair to assume that the errors of measurement for the different methods are uncorrelated.

In symbolic terms, δ_1 is error in measuring the independent variables by Method One, and ϵ_1 is error in measuring the dependent variables by Method One. $Var(\delta_1) = \Theta_{\delta_1}$ need not be diagonal, so errors of measurement for the independent variables may be correlated with one another. Similarly, $Var(\epsilon_1) = \Theta_{\epsilon_1}$ need not be diagonal, so errors of measurement for the dependent variables may be correlated with one another. And, errors of measurement using the same method may be correlated between the independent and dependent variables. For method one, this is represented by the matrix $Cov(\delta_1, \epsilon_1) = \Theta_{\delta\epsilon_1}$. The same pattern holds for Method Two. On the other hand, δ_1 and ϵ_1 are each independent of both δ_2 and ϵ_2 .

To emphasize an important practical point, the matrices Θ_{δ_1} and Θ_{δ_2} must be of the same dimension, and Θ_{ϵ_1} and Θ_{ϵ_2} must be of the same dimension – but none of the corresponding elements need be equal. This means that measurements of a variable by two different methods need not be equally precise.

Proof of model identification The following is typical of easier proofs for structural equation models. The goal is to solve for the model parameters in terms of elements of the variance-covariance matrix of the observable data. This shows the parameters are functions of the distribution, so that no two distinct parameter values could yield the same distribution of the observed data.

Collecting \mathbf{X}_1 , \mathbf{X}_2 , \mathbf{Y}_1 and \mathbf{Y}_2 into a single long data vector \mathbf{D} , we write its variance-covariance matrix as a partitioned matrix:

$$\Sigma = \left[\begin{array}{c|c|c|c} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} & \Sigma_{14} \\ \hline & \Sigma_{22} & \Sigma_{23} & \Sigma_{24} \\ \hline & & \Sigma_{33} & \Sigma_{34} \\ \hline & & & \Sigma_{44} \end{array} \right],$$

where the covariance matrix of \mathbf{X}_1 is Σ_{11} , the covariance matrix of \mathbf{X}_2 is Σ_{22} , the matrix of covariances between \mathbf{X}_1 and \mathbf{Y} is Σ_{13} , and so on.

Now we express all the Σ_{ij} sub-matrices in terms of the parameter matrices of Model (9) by straightforward variance-covariance calculations. Students may be reminded that things go smoothly if one substitutes for everything in terms of exogenous variables and

error terms before actually starting to calculate covariances. For example,

$$\begin{aligned}
\Sigma_{13} &= Cov(\mathbf{X}_1, \mathbf{Y}_1) \\
&= E(\mathbf{X}_1 \mathbf{Y}_1') \\
&= E((\boldsymbol{\xi} + \boldsymbol{\delta}_1)(\boldsymbol{\eta} + \boldsymbol{\epsilon}_1)') \\
&= E((\boldsymbol{\xi} + \boldsymbol{\delta}_1)(\boldsymbol{\Gamma}\boldsymbol{\xi} + \boldsymbol{\zeta} + \boldsymbol{\epsilon}_1)') \\
&= E((\boldsymbol{\xi} + \boldsymbol{\delta}_1)(\boldsymbol{\xi}'\boldsymbol{\Gamma}' + \boldsymbol{\zeta}' + \boldsymbol{\epsilon}_1')) \\
&= E(\boldsymbol{\xi}\boldsymbol{\xi}'\boldsymbol{\Gamma}' + \boldsymbol{\xi}\boldsymbol{\zeta}' + \boldsymbol{\xi}\boldsymbol{\epsilon}_1' + \boldsymbol{\delta}_1\boldsymbol{\xi}'\boldsymbol{\Gamma}' + \boldsymbol{\delta}_1\boldsymbol{\zeta}' + \boldsymbol{\delta}_1\boldsymbol{\epsilon}_1') \\
&= E(\boldsymbol{\xi}\boldsymbol{\xi}')\boldsymbol{\Gamma}' + E(\boldsymbol{\xi})E(\boldsymbol{\zeta}') + E(\boldsymbol{\xi})E(\boldsymbol{\epsilon}_1') + E(\boldsymbol{\delta}_1)E(\boldsymbol{\xi}')\boldsymbol{\Gamma}' + E(\boldsymbol{\delta}_1)E(\boldsymbol{\zeta}') + E(\boldsymbol{\delta}_1\boldsymbol{\epsilon}_1') \\
&= \boldsymbol{\Phi}\boldsymbol{\Gamma}' + 0 + 0 + 0 + 0 + \boldsymbol{\Theta}_{\delta\epsilon_1}.
\end{aligned}$$

In this manner, we obtain the partitioned covariance matrix of the observable data $\mathbf{D} = (\mathbf{X}'_1, \mathbf{X}'_2, \mathbf{Y}'_1, \mathbf{Y}'_2)'$ as

$$\begin{aligned}
\Sigma &= \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} & \Sigma_{14} \\ & \Sigma_{22} & \Sigma_{23} & \Sigma_{24} \\ & & \Sigma_{33} & \Sigma_{34} \\ & & & \Sigma_{44} \end{bmatrix} \\
&= \begin{bmatrix} \boldsymbol{\Phi} + \boldsymbol{\Theta}_{\delta_1} & \boldsymbol{\Phi} & \boldsymbol{\Phi}\boldsymbol{\Gamma}' + \boldsymbol{\Theta}_{\delta\epsilon_1} & \boldsymbol{\Phi}\boldsymbol{\Gamma}' \\ & \boldsymbol{\Phi} + \boldsymbol{\Theta}_{\delta_2} & \boldsymbol{\Phi}\boldsymbol{\Gamma}' & \boldsymbol{\Phi}\boldsymbol{\Gamma}' + \boldsymbol{\Theta}_{\delta\epsilon_2} \\ & & \boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}' + \boldsymbol{\Psi} + \boldsymbol{\Theta}_{\epsilon_1} & \boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}' + \boldsymbol{\Psi} \\ & & & \boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}' + \boldsymbol{\Psi} + \boldsymbol{\Theta}_{\epsilon_2} \end{bmatrix}
\end{aligned} \tag{10}$$

The equality (10) corresponds to a system of ten matrix equations in nine matrix unknowns. The unknowns are the parameter matrices of Model (9): $\boldsymbol{\Phi}$, $\boldsymbol{\Theta}_{\delta_1}$, $\boldsymbol{\Gamma}$, $\boldsymbol{\Theta}_{\delta\epsilon_1}$, $\boldsymbol{\Theta}_{\delta_2}$, $\boldsymbol{\Theta}_{\delta\epsilon_2}$, $\boldsymbol{\Psi}$, $\boldsymbol{\Theta}_{\epsilon_1}$ and $\boldsymbol{\Theta}_{\epsilon_2}$. Solving for them involves only simple substitutions, as follows.

$$\begin{aligned}
\boldsymbol{\Phi} &= \Sigma_{12} \\
\boldsymbol{\Gamma} &= \Sigma'_{13}\boldsymbol{\Phi}^{-1} = \Sigma'_{23}\boldsymbol{\Phi}^{-1} \\
\boldsymbol{\Theta}_{\delta_1} &= \Sigma_{11} - \boldsymbol{\Phi} \\
\boldsymbol{\Theta}_{\delta_2} &= \Sigma_{11} - \boldsymbol{\Phi} \\
\boldsymbol{\Theta}_{\delta\epsilon_1} &= \Sigma_{13} - \boldsymbol{\Phi}\boldsymbol{\Gamma}' \\
\boldsymbol{\Theta}_{\delta\epsilon_2} &= \Sigma_{24} - \boldsymbol{\Phi}\boldsymbol{\Gamma}' \\
\boldsymbol{\Psi} &= \Sigma_{34} - \boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}' \\
\boldsymbol{\Theta}_{\epsilon_1} &= \Sigma_{33} - \boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}' - \boldsymbol{\Psi} \\
\boldsymbol{\Theta}_{\epsilon_2} &= \Sigma_{44} - \boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}' - \boldsymbol{\Psi}
\end{aligned} \tag{11}$$

This shows that Model (9) is identified, so that if data are collected following the double measurement recipe, then the data analyst may proceed without giving further thought to model identification.

A notable feature of the identifying equations in (10) is that $\Sigma_{23} = \Sigma_{14}$. This constraint on the covariance matrix Σ arises from the model, and is called an *over-identifying restriction*. Such constraints arise whenever there are more identifying equations than unknowns. Even non-identified models may imply constraints — testable constraints — on the covariance matrix. This is an interesting side-issue we shall not pursue here. At any rate, an identified model with more identifying equations than unknowns is called *over-identified*.

Notice how the over-identifying restriction led us to write two expressions for Γ in the second line of (11). This was not necessary for the proof, but it allows us to clear up a point of confusion that sometimes arises. In an over-identified model, the parameter vector may be written in terms of Σ in more than one way. However, this does *not* mean that there are multiple solutions to the identifying equations, which would make the model non-identified. Multiple solutions means that more than one set of *parameter values* would satisfy the equations; this is not the case here.

4 Intercepts

We now expand Model (9) to include intercepts and non-zero expected values. We will see that this leads to complications that are seldom worth the trouble, and the classical models with zero expected value and no intercepts are usually preferable. Let

$$\begin{aligned}
 \eta &= \alpha + \Gamma\xi + \zeta \\
 \mathbf{X}_1 &= \nu_1 + \xi + \delta_1 \\
 \mathbf{X}_2 &= \nu_2 + \xi + \delta_2, \\
 \mathbf{Y}_1 &= \nu_3 + \eta + \epsilon_1 \\
 \mathbf{Y}_2 &= \nu_4 + \eta + \epsilon_2,
 \end{aligned}
 \tag{12}$$

where α , ν_1 , ν_2 , ν_3 and ν_4 are vectors of constants, and $E(\xi) = \kappa$. Everything else is as in Model (9). We call the terms $\nu_1 \dots, \nu_4$ *measurement bias*. For example, if one of the elements of \mathbf{X}_1 is reported amount of exercise, the corresponding element of ν_1 would be the average amount by which people exaggerate how much they exercise.

Again, the observable data \mathbf{X}_1 , \mathbf{X}_2 , \mathbf{Y}_1 and \mathbf{Y}_2 are collected into a data vector \mathbf{D} , with expected value μ and covariance matrix Σ . The pair (μ, Σ) is a function of the probability distribution of \mathbf{D} . If the parameter matrices of Model (12) are functions of μ and Σ , then they are also functions of the distribution of \mathbf{D} , and they will be identified.

Since the addition of constants has no effect on variances or covariances, the contents of Σ are given by (10), as before. The expected value $\boldsymbol{\mu}$ is the partitioned vector

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \\ \boldsymbol{\mu}_3 \\ \boldsymbol{\mu}_4 \end{bmatrix} = \begin{bmatrix} \frac{E(\mathbf{X}_1)}{E(\mathbf{X}_2)} \\ \frac{E(\mathbf{Y}_1)}{E(\mathbf{Y}_2)} \end{bmatrix} = \begin{bmatrix} \frac{\boldsymbol{\nu}_1 + \boldsymbol{\kappa}}{\boldsymbol{\nu}_2 + \boldsymbol{\kappa}} \\ \frac{\boldsymbol{\nu}_3 + \boldsymbol{\alpha} + \boldsymbol{\Gamma}\boldsymbol{\kappa}}{\boldsymbol{\nu}_4 + \boldsymbol{\alpha} + \boldsymbol{\Gamma}\boldsymbol{\kappa}} \end{bmatrix}. \quad (13)$$

To demonstrate the identification of Model (12), one would need to solve the equations in (13) uniquely for $\boldsymbol{\nu}_1$, $\boldsymbol{\nu}_2$, $\boldsymbol{\kappa}$ and $\boldsymbol{\alpha}$. Even with $\boldsymbol{\Gamma}$ considered known and fixed because it is identified in (10), this is impossible, because (13) specifies $2m + 2p$ additional equations in $3m + 3p$ additional unknowns.

It is tempting to assume the measurement bias terms $\boldsymbol{\nu}_1 \dots, \boldsymbol{\nu}_4$ to be zero; this would allow identification of $\boldsymbol{\alpha}$ and $\boldsymbol{\kappa}$. Unfortunately, we doubt that such an assumption could be justified very often in practice. And assuming no measurement bias when it is really present will have an unfortunate effect on the estimation of $\boldsymbol{\Gamma}$.

Most of the time, all we can do is identify the parameter matrices that appear in the covariance matrix, and also the *functions* $\boldsymbol{\mu}_1 \dots, \boldsymbol{\mu}_4$ of the parameters of Model (12). This can be viewed as a re-parameterization of the model.

Looking at (13), we see that the mean vector contains information about the regression coefficients in $\boldsymbol{\Gamma}$, but we can't get at it if there is measurement bias. We believe that except in special circumstances, this makes it reasonable to employ the classical no-intercept structural equation models to do regression with latent variables.

5 Normality

5.1 Maximum Likelihood Estimation

In the preceding section, we escaped from a serious model identification problem by re-parameterizing the means and intercepts. It is instructive to see how this works in the multivariate normal case, where the parameters would be estimated by maximum likelihood. For $i = 1, \dots, n$, we collect the observed data $\mathbf{x}_{i,1}$, $\mathbf{x}_{i,2}$, $\mathbf{y}_{i,1}$ and $\mathbf{y}_{i,2}$ into a vector \mathbf{d}_i , of length $2m + 2p$. We then write -2 times the log likelihood as a function of $\boldsymbol{\mu}$ and Σ .

$$\begin{aligned} -2 \log L(\boldsymbol{\mu}, \Sigma) &= -2 \log \prod_{i=1}^n \frac{1}{|\Sigma|^{\frac{1}{2}} (2\pi)^{\frac{k}{2}}} \exp \left[-\frac{1}{2} (\mathbf{d}_i - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{d}_i - \boldsymbol{\mu}) \right] \\ &= n [(2m + 2p) (\log |\Sigma| + \log 2\pi) + \text{tr}(\Sigma^{-1} \widehat{\Sigma}) \\ &\quad + (\bar{\mathbf{d}} - \boldsymbol{\mu})' \Sigma^{-1} (\bar{\mathbf{d}} - \boldsymbol{\mu})], \end{aligned} \quad (14)$$

where tr denotes the trace of a matrix: the sum of its diagonal elements.

The goal, of course, is to minimize (14) over all the model parameters making up $\boldsymbol{\mu}$ and Σ . Now for any value of Σ (so long as it is non-singular and hence positive definite),

the quadratic form in the last line of (14) is zero and the entire function is minimized when $\boldsymbol{\mu}$ equals $\bar{\mathbf{d}}$. This means that “centering the data” by subtracting off sample means and then pretending that all variables have expected value zero is equivalent to starting with a model like (12) that contains intercepts, re-parameterizing the components of $\boldsymbol{\mu}$ in (13) as $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_4$, and then estimating those functions by maximum likelihood.

Once this is done, the last line of (14) disappears, and inference is based solely on the sample covariance matrix. We are back to the no-intercept Model (9). Denoting the vector of model parameters by θ , the MLE is obtained by minimizing the $-2 \log$ likelihood

$$n[(2m + 2p)(\log |\boldsymbol{\Sigma}(\theta)| + \log 2\pi) + \text{tr}(\boldsymbol{\Sigma}(\theta)^{-1} \widehat{\boldsymbol{\Sigma}})]$$

numerically over all θ in the parameter space.

The full range of large-sample likelihood methods is then available. Maximum likelihood estimates are asymptotically normal, and asymptotic standard errors are convenient by-products of the numerical minimization; most software produces them by default. Dividing an estimated regression coefficient by its standard error gives a Z -test for whether the coefficient is different from zero. We should mention that our experience is that likelihood ratio tests can substantially outperform both these Z -tests and the Wald tests that are their generalizations, especially when there is a lot of measurement error and the variables are strongly related to one another.

5.2 Robustness of the Normal Model

In presenting models for regression with measurement error, it is sometimes convenient to assume that everything is multivariate normal. This is especially true when giving examples of models that are not identified. But normality is not necessary. Suppose Model (9) holds, and that the distributions of the latent independent variables and error terms are unknown, except for possessing covariance matrices. In this case the parameter of the model could be expressed as $\theta = (\boldsymbol{\Gamma}, \boldsymbol{\Phi}, \boldsymbol{\Psi}, \boldsymbol{\Theta}_{\delta_1}, \boldsymbol{\Theta}_{\delta_2}, \boldsymbol{\Theta}_{\epsilon_1}, \boldsymbol{\Theta}_{\epsilon_2}, \boldsymbol{\Theta}_{\delta\epsilon_1}, \boldsymbol{\Theta}_{\delta\epsilon_2}, F_{\boldsymbol{\xi}}, F_{\boldsymbol{\zeta}}, F_{\boldsymbol{\delta}_1}, F_{\boldsymbol{\delta}_2}, F_{\boldsymbol{\epsilon}_1}, F_{\boldsymbol{\epsilon}_2})$, where $F_{\boldsymbol{\xi}}, F_{\boldsymbol{\zeta}}, F_{\boldsymbol{\delta}_1}, F_{\boldsymbol{\delta}_2}, F_{\boldsymbol{\epsilon}_1}$ and $F_{\boldsymbol{\epsilon}_2}$ are the (joint) cumulative distribution functions of $\boldsymbol{\xi}, \boldsymbol{\zeta}, \boldsymbol{\delta}_1, \boldsymbol{\delta}_2, \boldsymbol{\epsilon}_1$ and $\boldsymbol{\epsilon}_2$ respectively.

Note that the parameter in this “non-parametric” problem is of infinite dimension, but this presents no conceptual difficulty. The probability distribution of the observed data is still a function of the parameter vector, and to show model identification, we would have to be able to recover the parameter vector from the probability distribution of the data. While in general we cannot recover the whole thing, we certainly can recover a useful *function* of the parameter vector, namely $\boldsymbol{\Gamma}$. In fact, $\boldsymbol{\Gamma}$ is the only quantity of interest; the remainder of the parameter vector consists only of nuisance parameters, whether the model is normal or not.

To make the reasoning explicit, the covariance matrix $\boldsymbol{\Sigma}$ is a function of the probability distribution of the observed data, whether that probability distribution is known or not. The calculations leading to (11) still hold, showing that $\boldsymbol{\Gamma}$ is a function of $\boldsymbol{\Sigma}$, and hence of the probability distribution of the data. Therefore, $\boldsymbol{\Gamma}$ is identified.

This is all very well, but can we actually *do* anything without knowing what the distributions are? Certainly! Looking at (11), one is tempted to just put hats on everything to obtain Method-of-Moments estimators. However, we can do a little better. Note that while $\Phi = \Sigma_{12}$ is a symmetric matrix in the population and $\widehat{\Sigma}_{12}$ *converges* to a symmetric matrix, $\widehat{\Sigma}_{12}$ will be non-symmetric for any finite sample size (with probability one). A better estimator is obtained by averaging pairs of off-diagonal elements:

$$\widehat{\Phi}_M = \frac{1}{2}(\widehat{\Sigma}_{12} + \widehat{\Sigma}'_{12}),$$

where the subscript M indicates a Method-of-Moments estimator. It is consistent by the Law of Large Numbers. Using the second line of (11), a reasonable though non-standard estimator of Γ is

$$\widehat{\Gamma}_M = \frac{1}{2}(\widehat{\Sigma}'_{13} + \widehat{\Sigma}'_{23})\widehat{\Phi}_M^{-1} \quad (15)$$

Consistency follows from the Law of Large Numbers and a continuity argument. All this assumes the existence only of second moments and cross-moments. With the assumption of fourth moments, the multivariate Central Limit Theorem would provide a routine basis for large-sample interval estimation and testing.

However, there is no need to bother. Research on the robustness of the normal model for structural equation model (Amemiya, Fuller and Pantula, 1987; Anderson and Rubin, 1956; Anderson and Amemiya, 1988; Anderson, 1989; Anderson and Amemiya, 1990; Browne, 1988; Browne and Shapiro, 1988; Satorra and Bentler, 1990) shows that procedures for (such as likelihood ratio and Wald tests) based on a multivariate normal model are asymptotically valid even when the normal assumption is false. And Satorra and Bentler (1990) describe Monte Carlo work suggesting that normal-theory methods generally perform better than at least one method (Browne, 1984) that is specifically designed to be distribution-free. Since the methods suggested by the estimator (15) are similar to Browne's weighted least squares approach, they are also likely to be inferior to the standard normal-theory tools.

It is important to note that while the normal-theory tests and confidence intervals for Γ can be trusted when the data are not normal, this does not extend to the other model parameters. For example, if the vector of latent variables ξ is not normal, then normal-theory inference about its covariance matrix will be flawed.

In any event, the method of choice is maximum likelihood, with interpretive focus on the regression coefficients in Γ rather than on the other model parameters.

6 Discussion

In general, data collection should be planned with the statistical analysis in mind. In keeping with this idea, the double measurement design is both a statistical model – specifically, Model (9) – and a set of guidelines for data collection. It assumes that

measurement error is present, and that when data are collected by a common method or in a common setting, the errors of measurement will naturally be correlated with one another. It also assumes that each independent variable can be measured more than once, ideally on different occasions and in ways that are different enough so that errors of measurement are independent between occasions.

A great deal of effort can be saved by following this recipe. The data are tailored to satisfy the technical requirements of the model, while the model allows for the inevitable correlations among measurement errors within occasions and is automatically identified, allowing clear conclusions to be drawn from the data. The only remaining issue is choosing good software and making sure that one knows what it is actually doing.

Unfortunately, most observational data sets are assembled without any awareness of measurement error as a statistical issue. Variables tend to be measured in only one way, and often at more or less the same time by the same personnel. Only after the data are collected do the investigators possibly start to think about fitting a model with measurement error. Many times, it is only at this point that a statistician enters the picture.

This is a difficult situation, but not necessarily hopeless. The most plausible model that includes measurement error is unlikely to be identified, but the instrumental variables example of Section 2 tells us that model identification can sometimes be purchased by adding more dependent variables. (Watch out, though! Dependent variables are usually measured with error too, and one needs specific reason to believe that those measurement errors are unrelated to measurement errors in the independent variables.) Sometimes, a model can be simplified or constrained, perhaps by assuming that certain covariances are zero, and the simplified model will be identified and still fairly realistic.

Fixing up a non-identified model after the data are already collected requires the quantitative sophistication to check model identification (repeatedly), and the subject-matter sophistication to tell whether the model is still scientifically meaningful when a given technical constraint is imposed. Either one person has to know a lot, or statistician and scientist must work closely together for an extended period, without any guarantee of ultimate success. It's a lot easier to plan the study properly in the first place.

One final comment is that from the statistician's viewpoint, a non-identified model is a "bad" model because it does not allow us to find out about the model parameters, and will probably generate a pile of warnings and error messages if we try to run the software anyway. But it's not the model's fault! Think of the very first example, the simple regression through the origin of Model (1). A model like this could be reasonable and even approximately correct, but the *data* we have will not allow us to estimate the parameters.

Now consider what happens when a fairly complicated initial model turns out not to be identified. The typical approach is to start imposing constraints that will make it identified. But this makes the model better only in a formal, statistical sense. Actually, the initial model was probably the most natural and believable one, and what we are doing is to chop pieces off for purely technical reasons. The best we can hope is that this

does not cripple the model too much.

It's not the model's fault; it's the data's fault. Or, to put it delicately, there is an opportunity for scientists to make their research even better by collecting data that allow reasonable models to be estimated.

References

Amemiya, Y., Fuller, W.A. and Pantula, S.C. (1987). "The asymptotic distributions of some estimators for a factor analysis model." *Journal of Multivariate Analysis*, 22, 51-64.

Anderson, T.W. (1989). "Linear latent variable models and covariance structures." *Journal of Econometrics*, 41, 91-119.

Anderson, T. W. and Amemiya, Y. (1988). "The asymptotic normal distribution of estimators in factor analysis under general conditions." *Annals of Statistics*, 16, 759-771.

Anderson, T. W. and Amemiya, Y. (1990). "Asymptotic chi-square tests for a large class of factor analysis models." *Annals of Statistics*, 18, 1453-1463.

Anderson, T. W. and Rubin, H. (1956). "Statistical inference in factor analysis." *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, 5, 111-150.

Arbuckle, J. L. (2006). *Amos 7.0 User's Guide*. Chicago: SPSS Inc.

Bentler, P. M. and Weeks, D. G. (1980). "Linear structural equations with latent variables." *Psychometrika*, 45, 289-308.

Bentler, P. M. and Woodward, J. A. (1978). "A Head Start re-evaluation: Positive effects are not yet demonstrable." *Evaluation Quarterly*, 2, 493-510.

Bollen, K. A. (1989). *Structural equations with latent variables*. New York: Wiley.

Browne, M.W. (1984). "Asymptotically distribution-free methods for the analysis of covariance structures." *British Journal of Mathematical and Statistical Psychology*, 37, 62-83.

Browne, M.W. (1987). "Robustness of statistical inference in factor analysis and related models." *Biometrika*, 74, 375-384.

Browne, M.W. and Shapiro, A. (1988). "Robustness of normal theory methods in the analysis of linear latent variable models." *British Journal of Mathematical and Statistical Psychology*, 41, 193-208.

Cheng, C. L. and Van Ness, J. W. (1999). *Statistical regression with measurement error*.

London: Chapman & Hall.

Cochran, W. G. (1968). "Errors of measurement in statistics." *Technometrics*, 10, 637-666.

Fisher, R. A. F. (1938). *Statistical methods for research workers (7th ed.)*. London: Oliver and Boyd.

Fox, J. (2006). "Structural equation modeling with the `sem` package in R." *Structural equation modelling*, 13, 465-486.

Fuller, W. A. (1987). *Measurement error models*. New York: Wiley.

Goldberger, A. S. and Duncan, O. D. (1973). *Structural equation models in the social sciences*. New York: Academic Press.

Jöreskog, K. G. (1978). "Structural analysis of covariance and correlation matrices." *Psychometrika*, 43, 443-477.

Jöreskog, K. G. and Sörbom, D. (1996). *LISREL 8: Structural equation modelling with the SIMPLIS command language*. London: Scientific Software International.

McArdale, J. J. (1980). "Causal modelling applied to psychonomic systems simulation." *Behavior research methods & Instrumentation*. 12, 193-209.

McDonald, R. P. (1978). "A simple comprehensive model for the analysis of covariance structures." *British Journal of Mathematical and Statistical Psychology*. 31, 59-72.

SAS Institute, Inc. (1999). "SAS/STAT User's guide, Version 8." Cary, N. C.: SAS Institute, Inc. 3884 pp.

Satorra, A. and Bentler, P. M. (1990). "Model Conditions for Asymptotic Robustness in the Analysis of Linear Relations." *Computational Statistics and Data Analysis*. 10, 235-249.

Stouffer, S. A. (1936). "Evaluating the effect of inadequately measured variables in partial correlation analysis." *Journal of the American Statistical Association*, 31, 348-360.