INDIAN INSTITUTE OF MANAGEMENT
AHMEDABAD • INDIA

# A Conceptual Overview of Structural Equation Modeling

**Tathagata Banerjee**
**Arindam Banerjee**
**Erina Paul**

विद्याविनियोगाद्विकासः

INDIAN INSTITUTE OF MANAGEMENT
AHMEDABAD-380 015
INDIA

# A Conceptual Overview of Structural Equation Modeling

**Tathagata Banerjee[1]**
**Arindam Banerjee[2]**
**Erina Paul[3]**

## Abstract

*A synthesized version of Structural Equation Modelling (SEM) and its possible applications in Management problems is presented. The main contribution of the paper is its simple description of a somewhat complex statistical process for the understanding of the beginners in this domain. It acts as a initial reading in SEM, before the researchers delve into more complex exposition of the statistical technique. The description is largely in English (not statistics) and is palatable to readers not trained enough in the domain of statistics.*

*It will serve as a good overview of this methodology for FPM students in business schools.*

---

[1] Professor, Indian Institute of Management, Ahmedabad  Email: tathagata@iimahd.ernet.in

[2] Professor, Indian Institute of Management, Ahmedabad  Email: arindam@iimahd.ernet.in

[3]  Erina Paul was a Research Assistant at Indian Institute of Management, Ahmedabad

## 1. Introduction:

Structural equation models, also called simultaneous equation models, refer to multiequation systems that include continuous **latent variables**[1] each representing a **concept or construct**[2], multiple **indicators**[3] of a concept or construct that may be continuous, ordinal, dichotomous or censored, errors of measurement and errors in equations. One may also view it as an interrelated system of regression equations where some of the variables (latent or observable) have multiple indicators and where measurement error is taken into account when estimating relationships. From a different point of view, these are factor analysis models in which factor loadings are restricted to zero or some other constants, and the researcher allows factors to influence each other, directly and indirectly. The most general form of the structural equation model includes Analysis of Variance, Analysis of Covariance, Multiple Linear Regression, Multivariate Multiple Regression, Recursive and Non-recursive Simultaneous Equations, Path Analysis, Confirmatory Factor Analysis and many other procedures as special cases. So, the term "Structural Equation Model" (SEM) refers to a comprehensive statistical methodology for testing and estimating causal relations using a combination of cross-sectional statistical data and qualitative causal assumptions. Unlike the usual multivariate linear regression model, the response variable in one regression equation in a structural equation model may appear as a predictor in another equation. Indeed, variables in a structural equation model may influence one-another reciprocally, either directly or through other variables.

Structural equation models have been discussed extensively in psychological science (Rabe-Hesketh e. al., 2004; Cole and Maxwell, 2003; Muthén, B., 1984; Bentler and Weeks, 1980; Bentler and Tanaka, 1982; Bentler and Freeman, 1983; Anderson and Gerbing, 1987, 1991)**,** econometrics (Krishnakumar and Nagar, 2008; Muthen, 1983), social sciences and quantitative behavioral sciences (Anderson, 1987; Muthen, 1982, 2002; Krishnakumar, 2007, 2008; Netemeyer and Bentler, 2001; Bauer, 2003) and management science (Gerbing and Anderson, 1984, 1988; Anderson, et. al., 1987; Anderson, 1987; Bagozzi, 1981; Fornell and Larker, 1981a, 1981b; Bagozzi and Fornell, 1989).

Unfortunately, however, researchers in many other areas of potential applications are relatively unfamiliar with the concept and its implementation. A more generous explanation for this is, SEM's are close to the kind of informal thinking about causal relation that is common in theorizing in psychological science, social science and management science and therefore, researchers in these areas find these models useful for translating such theories into data analysis.

In section 2 we begin with an example to illustrate the use of structural equation modeling and introduce path diagrams, which are essential tools for structural equation modeling. In Section 3 we introduce SEM's for the general case for both causal model and the measurement model. In this generalized setting, we discuss the model identification in section 4. In Section 5 we briefly review the estimation of model parameters. In section 6 we consider the model evaluation and indices of model fits of SEM to data. In Section 7 we consider the specification problem of measurement model and in this context briefly discuss the concepts of validity, reliability and unidimensionality. In Section 8 we present a brief but important

discussion on whether SEM can be considered as a causal model. We present concluding remarks in Section 9.

## 2. An Example:

We consider an example discussed in **Bollen (1989**) to illustrate the use of SEM in building a **theory**[4] from data. The theory we consider here is, "Industrialization in developing countries is thought to enhance the chances of political democracy". Here Industrialization and Political Democracy are two constructs and the above theory is hypothesizing a relationship between them. It is often called causal relationship. At this point we avoid getting into a debate on what would be the proper definitions of the constructs, and what would be the right proxies or indicators of it. These are, of course, important issues that need to be considered seriously by every researcher at the outset. It is a common experience of the researchers that after the preliminary data collection the theory often is not validated by the data. This may happen because of faulty definitions of the constructs and/or due to wrong choices of the indicators and/or due to wrong choice of the causal model. The researcher then needs to revise either the definition of the constructs and/or the choice of the indicators and/or the causal model itself. We will discuss these issues at the end. However, at present we assume that the constructs are properly defined, the indicators are correctly chosen and the causal model is correctly specified.

We define Industrialization as "the degree to which a society's economy is characterized by mechanized manufacturing process", and Political Democracy, as "the extent of political rights and political liberties in a country". Both these constructs are unobservable and are thus represented in our model by what are called latent variables or unobserved variables. Our problem is to build and then test the above theory. Suppose we consider three latent variables, Industrialization in 1960 ($\xi_1$), Political Democracy in 1960 ($\eta_1$) and Political Democracy in 1965 ($\eta_2$). Here Industrialization is an **exogenous latent variable**[5] and Political Democracy is an **endogenous latent variable**[6]**.** The latent endogenous variables are only partially explained by the model and the unexplained part, i.e., the random disturbance in the equation is represented by $\zeta_i$. We assume that $\eta_2$ is a function of both $\xi_1$ and $\eta_1$. Also $\eta_1$ is a function $\xi_1$. Thus we have two equations expressing the above causal relationships.

$$\eta_1 = \gamma_{11} \, \xi_1 + \zeta_1$$
$$\eta_2 = \beta_{21} \, \eta_1 + \gamma_{21} \, \xi_1 + \zeta_2 \tag{2.1}$$

where, $\gamma_{11}$ , $\beta_{21}$ , $\gamma_{21}$ are structural parameters and have usual interpretations as in regression analysis, $\zeta_1$ and $\zeta_2$ are random disturbances with mean zero and are uncorrelated with the exogenous variable $\xi_1$. The latter assumption is necessary to avoid omitted variable bias.

Notice that here the equations are linear in variables and linear in parameters. Non-linear models are not much in use. Also the variables are expressed as deviation from its mean values. So the intercept term is absent. In matrix notation (2.1) can be re-written as,

$$\begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ \beta_{21} & 0 \end{bmatrix} \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix} + \begin{bmatrix} \gamma_{11} \\ \gamma_{21} \end{bmatrix} [\xi_1] + \begin{bmatrix} \zeta_1 \\ \zeta_2 \end{bmatrix} \tag{2.2}$$

or equivalently,

$$\boldsymbol{\eta} = \mathbf{B}\boldsymbol{\eta} + \boldsymbol{\Gamma}\boldsymbol{\xi} + \boldsymbol{\zeta} \tag{2.3}$$

where $\mathbf{E}\ (\zeta) = \mathbf{0}$ and $\mathbf{E}(\zeta\ \xi^{\mathbf{T}}\ ) = \mathbf{0}$. The diagonal elements of $\mathbf{B}$ are zero and the matrix $(\mathbf{I - B})$ is non-singular. Model (2.3) is known as **causal model** in SEM.

Now, observed variables are not perfectly correlated with the latent variables that they measure unless the latent variables are themselves observable. Nearly all measures of abstract concepts have far from perfect association with the latent variables they represent. We thus have the following **measurement model** linking the latent variables with the observable variables or proxies or indicators. The relations are imperfect rather than deterministic ones.

Let us consider the following indicators for Industrialization & Political Democracy. Three indicators for Industrialization we consider are, GNP per capita $(x_1)$, Inanimate Energy Consumption per capita $(x_2)$, Percentage of labor force in industry $(x_3)$ and indicators for Political Democracy are, expert ratings of the freedom of the press ($y_1$ in 1960, $y_5$ in 1965), expert ratings of the freedom of political opposition ($y_2$ in 1960, $y_6$ in 1965), expert ratings of the fairness of election ($y_3$ in 1960, $y_7$ in 1965), expert ratings of the effectiveness of the elected legislature ($y_4$ in 1960, $y_8$ in 1965). So the specification of the measurement model is

$$x_1 = \lambda_1\ \xi_1 + \delta_1\ , x_2 = \lambda_2\ \xi_1\ + \delta_2\ , x_3 = \lambda_3\ \xi_1\ + \delta_3 \qquad (2.4)$$

$$y_1 = \lambda_4\ \eta_1 + \varepsilon_1,\ y_5 = \lambda_8\ \eta_2 + \varepsilon_5,$$

$$y_2 = \lambda_5\ \eta_1 + \varepsilon_2,\ y_6 = \lambda_9\ \eta_2 + \varepsilon_6,$$

$$y_3 = \lambda_6\ \eta_1 + \varepsilon_3,\ y_7 = \lambda_{10}\ \eta_2 + \varepsilon_7,$$

$$y_4 = \lambda_7\ \eta_1 + \varepsilon_4,\ y_8 = \lambda_{11}\ \eta_2 + \varepsilon_8. \qquad (2.5)$$

In the above measurement model, $x_1$ to $x_3$ stand for indicators of $\xi_1$, $y_1$ to $y_4$ are the indicators of $\eta_1$, $y_5$ to $y_8$ are indicators of $\eta_2$. The $\lambda_i$'s are regression coefficients of the latent variables on the observed variables. The $\delta_i$'s and $\varepsilon_i$'s are the errors of measurements for $x_i$ and $y_i$. So in matrix notation we can write (2.4) and (2.5) as

$$\mathbf{x} = \mathbf{\Lambda_x}\ \mathbf{\xi} + \mathbf{\delta},\ \mathbf{\Lambda_x} = \text{diag}\ (\lambda_1, \lambda_2, \lambda_3\ )$$

$$\mathbf{y} = \mathbf{\Lambda_y}\ \mathbf{\eta} + \mathbf{\varepsilon},\ \mathbf{\Lambda_y} = \text{diag}\ (\lambda_4, \ldots, \lambda_{11}\ )\mathbf{.} \qquad (2.6)$$

Now we depict the system of simultaneous system of equations given by (2.3) and (2.6) using a **path diagram**[7] shown in Figure 1. The **direct effect**[8] of Industrialization 1960 on Political Democracy 1965 is $\gamma_{21}$. The **indirect effect**[9] of Industrialization 1960 on Political Democracy 1965 is $\gamma_{11}\ \beta_{21}$**.** Thus the total effect of Industrialization 1960 on Political Democracy 1965 is $\gamma_{21} + \gamma_{11}\ \beta_{21}$. If we consider regression of Political Democracy 1965 on Industrialization 1960 without bringing in the intervening variable Political Democracy 1960 in the model then the estimate of direct effect estimates the overall effect rather than the direct effect $\gamma_{21}$. This might lead to incorrect inference since even if $\gamma_{21}\gamma_{11}$ and $\beta_{21}$ are significantly different from zero the total effect may not be significant. It is often called the effect of omitting intervening

variable. Identifying intervening variables is crucial in building the model. It demands intimate domain knowledge and hard thinking on the researcher's part.

### 3. General Structural Equation Model

SEM has two components viz., the causal model and the measurement model. The causal model shows the linear relation between the latent variables (or equivalently constructs) while the measurement model shows the relation between the indicators and the latent variables.

**Causal Model:**

In general the causal model is written as follows,

$$\boldsymbol{\eta} = \boldsymbol{\alpha} + \mathbf{B}\boldsymbol{\eta} + \boldsymbol{\Gamma}\boldsymbol{\xi} + \boldsymbol{\zeta} \tag{3.1}$$

where $\boldsymbol{\alpha}$ is a vector of intercept term, $\boldsymbol{\eta}$ is a m×1 vector of endogenous latent variables and $\boldsymbol{\xi}$ is an n×1 vector of exogenous latent variables with mean $\boldsymbol{\kappa}$ and covariance matrix $\boldsymbol{\Phi}$, $\boldsymbol{\Gamma}$ is the m x n coefficient matrix for the effects of $\boldsymbol{\xi}$ on $\boldsymbol{\eta}$ , $\mathbf{B}$ is the m x m coefficient matrix showing the influence of the latent endogenous variables on each other and $\boldsymbol{\zeta}$ is a m×1 vector of error terms that has zero mean and covariance matrix $\boldsymbol{\Psi}$, and $\mathbf{cov}(\boldsymbol{\xi}, \boldsymbol{\zeta}) = \mathbf{0}$ . Usually the latent variables $\boldsymbol{\xi}$ and $\boldsymbol{\eta}$ are assumed to be measured as deviation from its means hence $\boldsymbol{\alpha} = \mathbf{0}$ and the model reduces to

$$\boldsymbol{\eta} = \mathbf{B}\boldsymbol{\eta} + \boldsymbol{\Gamma}\boldsymbol{\xi} + \boldsymbol{\zeta}.^* \tag{3.2}$$

In classical econometrics the simultaneous equation model is given by:

$$\mathbf{B}\mathbf{y} + \boldsymbol{\Gamma}\mathbf{x} = \mathbf{u},$$

where $\mathbf{y}$ represents a vector of exogenous variables and $\mathbf{x}$ a vector of endogenous variables.
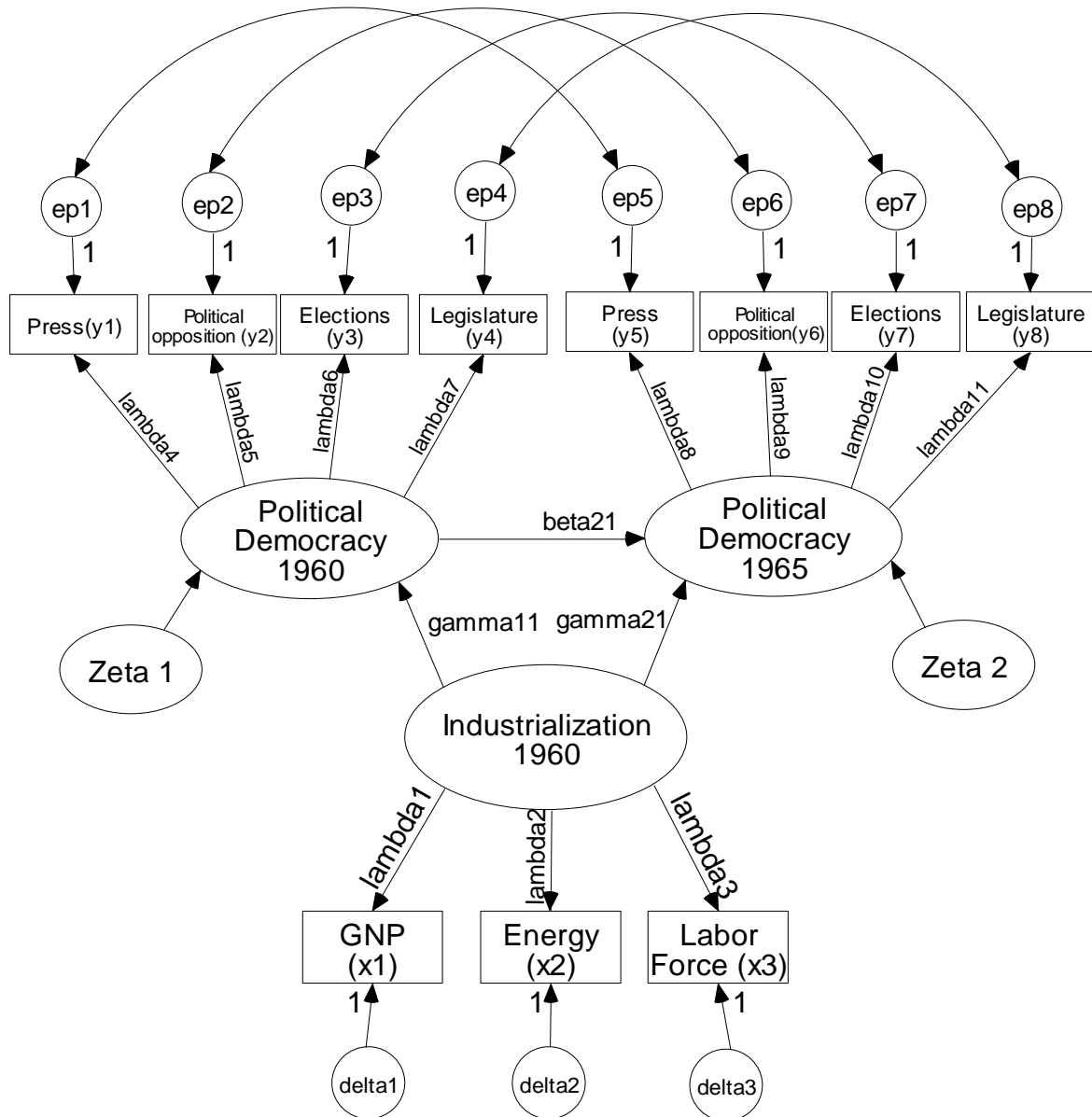
_____

[*]Since the model assumes $\mathbf{I} - \mathbf{B}$ is nonsingular, setting $\mathbf{A} = (\mathbf{I} - \mathbf{B})^{-1}$, it follows that

$$\boldsymbol{\eta} = \mathbf{A}\boldsymbol{\Gamma}\boldsymbol{\xi} + \mathbf{A}\boldsymbol{\zeta}, \quad \boldsymbol{\mu}_{\boldsymbol{\eta}} = \mathbf{A}\boldsymbol{\Gamma}\boldsymbol{\kappa}, \quad \mathbf{Cov}(\boldsymbol{\eta}) = \mathbf{A}(\boldsymbol{\Gamma}\boldsymbol{\Phi}\boldsymbol{\Gamma}^{\mathbf{T}} + \boldsymbol{\Psi})\mathbf{A}^{\mathbf{T}},$$

where $\boldsymbol{\mu}_{\boldsymbol{\eta}}$ is the mean vector and $\mathbf{Cov}(\boldsymbol{\eta})$ is the covariance matrix of $\boldsymbol{\eta}.$

# *Representation of the System of Simultaneous Equations by Path Diagram*

Note that the latent variable model (3.2) may be equivalently written as

$$\mathbf{A}\,\boldsymbol{\eta} + (\mathbf{-\Gamma})\,\boldsymbol{\xi} = \boldsymbol{\zeta}.$$

It is similar to the simultaneous equation model except that it is written in terms of latent or unobservable variables.

**Measurement Model:**

A test of the theory that the causal model formalizes is possible if we collect data on observable measures or indicators of the latent variables. The measurement model specifies the relation between the indicators and the latent variables. Suppose **y** represents the vector of p endogenous observed variables that are indicators of $\boldsymbol{\eta}$, and **x** represents the vector of q exogenous observed variables that are indicators of $\boldsymbol{\xi}$. We assume that these are expressed as deviation from its mean. The measurement model may then be expressed as

$$\mathbf{y} = \boldsymbol{\Lambda_y}\boldsymbol{\eta} + \boldsymbol{\varepsilon},$$
$$\mathbf{x} = \boldsymbol{\Lambda_x}\boldsymbol{\xi} + \boldsymbol{\delta}, \tag{3.4}$$

where $\boldsymbol{\Lambda_y}$ and $\boldsymbol{\Lambda_x}$ are the coefficient matrices, $\boldsymbol{\varepsilon}$ and $\boldsymbol{\delta}$ are the errors of measurements for **x** and **y**. Also $\boldsymbol{\varepsilon}$ and $\boldsymbol{\delta}$ are assumed to be uncorrelated with $\boldsymbol{\xi}$ and $\boldsymbol{\eta}$ and with each other. We further assume $\mathbf{E(\boldsymbol{\varepsilon})= E(\boldsymbol{\delta}) = E(\boldsymbol{\eta}) = E(\boldsymbol{\xi}) = 0, Cov(\boldsymbol{\varepsilon}) = \boldsymbol{\Theta_\varepsilon,} and\ Cov(\boldsymbol{\delta}) = \boldsymbol{\Theta_\delta}}$ where $\boldsymbol{\Theta_\varepsilon}$ and $\boldsymbol{\Theta_\delta}$ are the covariance matrices of $\boldsymbol{\varepsilon}$ and $\boldsymbol{\delta}$ respectively. If $\boldsymbol{\Theta_\varepsilon}$ and $\boldsymbol{\Theta_\delta}$ are equal to zero, the structural model (3.2) reduces to a simultaneous equation model.

After doing some algebra, one can express the covariance matrices of **y**, **x** and between **y** and **x** in terms of the unknown model parameters $\boldsymbol{\Lambda_x, \Lambda_y, \Phi, \Gamma, \Theta_\varepsilon}$ **and** $\boldsymbol{\Theta_\delta}$ which are given in Appendix 1. It may be worth reiterating that the empirical evidence of a possible relationship among the unobserved constructs is housed in $\Sigma$ (cf. Appendix 1), the covariance matrix of the observable (**y**, **x**). The sample covariance of (**y**, **x**) being a natural estimate of $\Sigma$ is then plugged into the relationships (A.1) (cf. Appendix 1) to estimate the value of each of the unknown model parameters and hence the necessity of deriving the components of the aforesaid matrix $\Sigma$ in relation to the unknowns.

It is worth reiterating that the distinctiveness of the SEM is the separation of the causal model from the measurement model. This is a broader generalization which does not constrain the model to assume unique and mono-dimensional measures for the constructs in the causal model. The flexibility available in the model construction has an appeal, since most research constructs (especially in psychological domain) are complex composites of elemental measures where the basis of their composition are not apparent up front. SEM does not require an a priori basis to construct measures of the latent variables, other than a classification of the measures to their respective latent constructs. This ensures that the value of the composite index of each latent construct is empirically computed (estimated) from the data. This flexibility comes at a cost to the researcher in terms of complications that may arise at the time of estimation. We shall discuss the issue in the following sections.

## 4. Model Identification

Following Hayduk (1987) let us illustrate the general issue of identification with an example not related to SEM. Suppose a theory claims that sum of two coefficients, say, $\alpha$ and $\beta$ must be equal to some specific number, while the available data indicate that the number is 10. Hence $\alpha + \beta = 10$, but what are the values of $\alpha$ and $\beta$? There are infinitely many such choices available that fit exactly. Thus the model constraints (requiring the sum of $\alpha$ and $\beta$) and the data constraints (the sum equaling 10) eliminate some sets of estimates of $\alpha$ and $\beta$ (for example 5-6, 7-9 etc.), but the combined effect of both constraints is insufficient to determine unique values for $\alpha$ and $\beta$. *It is the failure of the combined model and data constraints to identify (locate or determine) unique estimates that results in the name "the identification problem"* (Hayduk (1987)). In such a case the only option is to impose further model constraints or data constraints with the hope of eliminating more pairs of estimates and in the process eliminate the said (under) identification problem.

Such problems do crop up in the estimation of the SEM given that there are multiple unknown parameters and given the non-linear relationship among the variables, their interactions may yield innumerable permutations of possible estimate yielding the same outcome. An example of a similar kind is illustrated in the Appendix 2 to reinforce the notion of under identification that poses problems in the SEM.

Without claiming to be rigorous we may say that for an under-identified model estimate of at least one of the model parameters is not unique and thus unreliable; for an exactly identified model unique estimates for all the model parameters are available but estimates of its standard errors are not available; and for an over-identified model the estimates of the parameters and also the estimates of its standard errors are available. More the over identification the better it is. In order to build an over identified model starting from an under-identified one, either the number of model parameters be reduced by imposing meaningful constraints on the model, or the number of indicators or proxies for each latent variable be increased. For a detailed technical discussion on model identifiability problem we refer to Bollen (1989).

## 5. General Method of Estimation:

For a general structural equation model we always have an equation like (A.1) where on the left hand side we have $\Sigma$, the covariance matrix of the observed variables and on the right hand side we have a covariance matrix $\Sigma(\theta)$ involving the unknown vector of structural equation model parameters $\theta$. In (A.1) the unknown model parameters $\Lambda_x$, $\Lambda_y$, $\Phi$, $\Gamma$, $\Theta_\varepsilon$ **and** $\Theta_\delta$ comprise $\theta$. The matrix $\Sigma(\theta)$ is obtained by using the structural equations. So, in general we solve the equation

$$\Sigma = \Sigma(\theta) \tag{5.1}$$

for $\theta$ after replacing $\Sigma$ by its sample counterpart S. Here S is the sample covariance matrix of the observed variables.

For an over identified model we have more than one estimate of $\theta$ satisfying equation (5.1). So the question that naturally arises is which one should we choose? An obvious answer is the estimate that in some sense would minimize the discrepancy between **S** and $\Sigma$ ($\theta$).

In layperson's terms, imagine a search algorithm (with some intelligence) which searches across innumerable possible values of $\theta$ and hence of $\Sigma$ ($\theta$), with the aforesaid intelligence, to identify the value of $\theta$, say, $\hat{\theta}$ that gives minimum discrepancy between S and $\Sigma$ ($\hat{\theta}$). A typical exposition of a measure of closeness (minimum discrepancy) between S and $\Sigma$ is provided in Appendix 3 for reference.

## 6. Model Evaluation and model fits

If the Structural Equation Model is true then $\Sigma = \Sigma$ ($\theta$). Ideally we would then expect **S** – $\Sigma$ ($\hat{\theta}$), the sample covariance residual matrix, to be approximately null. But this is not the case since sample residuals are affected by several factors:

1. The departure of $\Sigma$ from $\Sigma$ ($\theta$)

2. The scales of observable variables

3. Sampling fluctuations

We are interested primarily to detect whether the structural equation model is correct, or more precisely how good the model fits the data.

Denoting the $(i, j)$-th element of S and $\Sigma$ ($\hat{\theta}$) by $s_{ij}$ and $\hat{\sigma}_{ij}$ respectively, we obtain the $(i, j)$-th element of the residual matrix as $s_{ij}$ - $\hat{\sigma}_{ij}$. If all the residuals in the residual matrix are positive, the model plausibly underpredicts covariance; if all the residuals are negative then the model plausibly overpredicts covariance. If $\Sigma \neq \Sigma$ ($\theta$) then lack of such incongruity may manifest itself in the sample residuals.

As stated above the departure of $\Sigma$ from $\Sigma$ ($\theta$) may also arise due to differences in the scales of the observed variables. A larger sample residual may arise since the scale units of one may have a much larger range than the others. In fact, if the ranges of the observed variables are too different, it may distort the comparison of the residuals. A simple solution to this is to calculate the correlation residuals $r_{ij} - \hat{r}_{ij}$, where $r_{ij}$ and $\hat{r}_{ij}$ are $(i, j)$-th element of the correlation matrices obtained from S and $\Sigma$ ($\hat{\theta}$) respectively.

To take care of the effect of sample sizes on sampling fluctuations as well as of scales Jöreskog and Sörbom (1986) proposed a standardized residual for each component of the residual matrix. It is given by:

$$\frac{(S_{ij} - \sigma_{ij}(\hat{\theta}))}{[(\sigma_{ii}(\hat{\theta})\sigma_{jj}(\hat{\theta}) + \sigma_{ij}^2(\hat{\theta})) / N]^{1/2}}.$$

The numerator represents the residual and the denominator its approximate standard error. The largest numerical value of the standardized residual indicates the element that is given

the worst fit by the model. Jöreskog and Sörbom (1986) also propose Goodness of Fit Index (GFI) and Adjusted Goodness of Fit Index (AGFI) for indicating the overall fit of the model if the model is fitted using maximum likelihood method. Writing $\Sigma(\hat{\theta}) = \hat{\Sigma}$, these are given by:

$$GFI_{ML} = 1 - \frac{tr[(\hat{\Sigma}^{-1}S - I)^2]}{tr[(\hat{\Sigma}^{-1}S)^2]}$$

$$AGFI_{ML} = 1 - [\frac{q(q+1)}{2df}][1 - GFI_{ML}].$$

Note that AGFI is GFI corrected for degrees of freedom of the model. Similarly they propose indices for the models fitted with unweighted and weighted least squares methods. For perfect fit, i.e., when, $\hat{\Sigma} = S$ these indices are equal to unity. Tests of hypotheses can be carried out for testing $\Sigma = \Sigma(\theta)$ and also for testing a sequence of nested models. For a detailed discussion we refer to Bollen (1989).

## 7. Specification of the Measurement Model

The development of SEM with latent variables has provided researchers especially in the realm of social science with considerable means to build, test, and modify theories. Jöreskog (1970, 1978) used maximum likelihood method to estimate the parameters of measurement and causal models simultaneously from the observed correlation (or covariance) matrix. Later they implemented the methodology in LISREL (Jöreskog and Sörbom (1978)), a widely used software of SEM. However, the initial analysis almost invariably indicates the need for a revision of either the measurement model, or the causal model or both. It is always advisable to think of the modeling process as the analysis of two conceptually distinct models: measurement model and causal model (Gerbing (1979), Jöreskog and Sörbom (1978)). The reason for drawing a distinction between the measurement model and the causal model is that proper specification of the measurement model is necessary before meaning can be ascribed to the analysis of causal model.

The SEM provides the flexibility to construct a generalized model where, the measurement of the constructs is not tightly defined by the distinct measures (for instance, if unidimensionality is not prescribed). While conceptually such models can exist and, is to an extent a cause for the growing popularity among some segment of researchers, lack of distinct identity of the construct by its corresponding measures can lead to severe (under) identification problems during the estimation the model. Hence, a good measurement model of the latent variables is prerequisite to the analysis of causal relations among the latent variables.

There are four stages in the specification of a measurement model. First, a theoretical definition of each construct should be put forward. A theoretical definition explains in simple and precise terms the meaning of a construct. In the example introduced at the outset, the construct 'Political Democracy' is defined as "the extent of political rights and political liberties in a country." Once we define the construct, its dimensions are identifiable. Dimensions are the distinct aspects of a construct that is not further divisible into components. The dimensions attached with the construct 'Political Democracy' are, 'Political Rights' and 'Political Liberty'.

Now for each dimension in the example above two indicators or observable measures are considered. For 'Political Rights' the measures considered are 'Expert Ratings of the Fairness of Election' and 'Expert Ratings of the Effectiveness of the Elected Legislature' and for 'Political Liberty' these are 'Expert Ratings of the Freedom of Political Opposition' and 'Expert Ratings of the Fairness of Election'. Once the measures are identified, we specify the relation between the measures and the latent variables.

Thus the four steps in the specification of a measurement model are, (i) define the constructs, (ii) identify the dimensions and the latent variables, (iii) find the measures or indicators and, (iv) specify the relation between the measures and the latent variables.

In testing a theory using SEM, once the concepts or theoretical constructs are defined the researcher estimates each construct using a posited relation between it and multiple indicators. The estimation and testing of the posited relationship by using SEM methodology is often called the **confirmatory factor analysis**[10] (Holzinger, 1944; Jöreskog 1966, 1969) in contrast to **exploratory factor analysis**[11].

However, during the model specification, the researcher has to answer a few important questions. Most important is; do the selected indicators measure the construct they are supposed to measure? In other words, are these indicators valid measures of the underlying construct? Further, validity of the indicators is not enough to ensure a good specification of the measurement model. Even if these are valid, we need to verify whether the indicators are reliable. By reliability of an indicator we mean, if the indicator is measured repeatedly, the measurements should be consistent. So the specification of a measurement model may be reliable without being valid or may be valid without being reliable. For example, if a faulty instrument measures weight always five kg less than the actual, the measurements are not valid but reliable. On the other hand, if the instrument gives highly variable measurements centred on the actual weight, the measurements are valid but not reliable. During the measurement process viz., specification of the measurement model, it is thus important to verify whether the measurement process is valid and reliable. If a measurement process fails in ensuring either validity or reliability or both, the estimated causal relationships between the constructs would consequently be invalid or unreliable or both. Various concepts of validity and measures of reliability are available in the literature. For a detailed discussion we refer to Bollen (1989). Finally, the researcher needs to verify whether the set of indicators defining each construct is **unidimensional or congeneric**[12] (Aaker and Bagozzi (1979), Bagozzi (1980) p.125-8; Jöreskog (1970)). Lack of unidimensionality most often represents a measurement model misspecification and unfortunately, a number of misspecifications of this kind typically occur with initial models. There are various methods proposed (Anderson and Gerbing (1982), Anderson, Gerbing and Hunter (1987), Anderson and Gerbing (1991)) in the literature to verify unidimensionality of the measurement model. However, some researchers feel (Bagozzi and Fornell (1989)) unidimensionality is a concept difficult to establish empirically.

## 8. SEM and Causality

Following Bollen (1989) let us briefly discuss the concept of causality. Consider a variable y which is isolated from all influences except from a second variable called x. If a change in y

accompanies a change in x, then x is associated with y. In order to establish x causes y, we must ensure that the association is due to x causing y not the other way around. The definition of causality has thus three components, (i) Isolation, (ii) Association and (iii) Direction of Causation.

We consider a simple illustrative example. Suppose y represents the incidence of lung cancer for each state in India. An argument that may be put forward is, since every case of lung cancer has a unique and unpredictable origin, *y* is a random disturbance term representing the sum total effect of innumerable infinitesimal causes.
So the model is

$$y = \zeta \qquad (7.1)$$

where $\zeta$ represents the disturbance term. This model represents the position of an extreme skeptic who believes *y* is incapable of being systematically explained by other variables. On the other hand a closer look at $\zeta$ may lead to the discovery of one or more variables, which could be meaningfully brought into the model. The simplest assumption may be, $\zeta$ consists of a single variable, say, the number of smokers (x) in the state. We then assume a simple model *y* = *f*(*x*) connecting y to x. To make it simpler, we assume

$$y = \beta x. \qquad (7.2)$$

But most of us would feel uncomfortable with model (7.2). It looks like an assumption almost certainly not true. Most of us may be comfortable with a model which is between the two extremes represented by (7.1) and (7.2). A reasonable model is,

$$y = \beta x + \zeta. \qquad (7.3)$$

Note model (7.3) clearly violates the condition of isolation. The disturbance term $\zeta$ is unobserved. We cannot control it. Isolation being impossible, we define what is called *pseudo isolation* condition. To assume *x* is isolated from $\zeta$, a simple assumption may be *x* is uncorrelated with $\zeta$. This is the condition implicitly assumed by most classical econometric models. However, pseudo isolation is nearly impossible to attain because of left out intervening variables, reciprocal causation, wrong model specification, presence of measurement errors, correlated disturbances, nonrandom sample selection etc. Regarding the direction of causality, the single most effective means of proving it is to establish temporal priority. This does not always work. Also it is not always clear that temporal priority is met especially when the models involve latent variables and its indicators. To sum up, proving causality beyond any doubt does not seem to be a practical proposition.

A misconception that is prevalent among the SEM users is that it establishes causality. Our discussion above shows that for demonstrating causality isolation from the effects of other variables must be ensured, association must be demonstrated and direction of causality should be established. In almost all applications of SEM these conditions are not met. Most SEM applications are best viewed as potential explanations for whether the causal relationship

envisaged in the model is consistent with the data. Many researchers (Bullock et. al., 1994; Hoyle and Smith, 1994; MacCallum et.al., 1993) argued causal inferences from such models are rare and possibly ill advised. Through putative logic, strong theoretical arguments, and longitudinally collected data one can strengthen a causal argument. Ultimately it is the design, not the statistical method, (i.e., SEM), that permits causal hypotheses to be adequately tested (Bullock et. al., 1994; Campbell & Stanley, 1963; Hoyle and Smith, 1994; Sobel, 1993). Cliff (1983) presented a sobering reminder on this issue by stating, "data do not confirm a model, they only fail to disconfirm it". It indeed echoes Popper's view, verification is impossibility, only falsification is possible. It is indeed unfortunate that numerous articles are written from a perspective, as if, we seek to confirm that our models fit. Cliff continues further, "when data do not disconfirm a model, there are many other models that are not disconfirmed either". MacCullam et.al.(1993) demonstrated that there were astronomical number of alternative models published in prestigious journals that would have provided the equivalent fit. If we seek to make a causal statement, we would best operate experimentally. "The most satisfactory, almost the only satisfactory, method for demonstrating causality is the active control of variables" (Cliff, 1983).

## 9. Concluding Remarks

A researcher's model should pass the tests of both "Model-Data Consistency" and "Model-Real-world Consistency" (Bollen, 1989) in order to be relevant and useful. Unfortunately, however, most applications of SEM test the former and only implicitly assume the latter. The reason is, checking "Model-Data Consistency" is considered to be an inseparable part of SEM methodology and is checked by looking at discrepancy between $\hat{\Sigma}$ and S, its magnitude, sign and statistical significance. On the other hand "Model-Reality Consistency" is a more "slippery" issue and is not directly verifiable from data. Here the question that a researcher should ask is, does the model mirror the real world? Checking this consistency thus needs intimate knowledge of how the 'real' world works. In practice we imperfectly evaluate the "Model-Real-world Consistency" of a model by its predictive validity (its power in predicting future events) or by cross-validating (validating the model) it with independent data sets. "It is tempting to use model-data consistency as proof of model-reality consistency" (Bollen, 1989), but it would be misleading. "Model-Reality Consistency" clearly implies "Model-Data Consistency" but rarely the other way around.

## Appendix 1

## Covariance Matrix of the Observed Variables:

We provide the algebraic formulations of various components of the information matrix $\Sigma$.

$$\Sigma_{YY} = \Lambda_Y [A (\Gamma \Phi \Gamma^T + \Psi) A^T] \Lambda_Y^T + \Theta_\varepsilon$$

$$\Sigma_{XX} = \Lambda_X \Phi \Lambda_X^T + \Theta_\delta \qquad\qquad (A.1)$$

$$\Sigma_{YX} = \Lambda_Y A \Gamma \Phi \Lambda_X^T,$$

where cov $(\mathbf{y}) = \Sigma_{YY}$, cov$(\mathbf{x}) = \Sigma_{XX}$ and cov $(\mathbf{y}, \mathbf{x}) = \Sigma_{YX}$. From (3.4) - (3.7), it follows that the covariance matrix of the observed variables may be expressed as,

$$\Sigma = \text{cov}(y, x) = \begin{bmatrix} \Sigma_{YY} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{XX} \end{bmatrix}.$$

An important special case where $\mathbf{y}$ and $\mathbf{x}$ are observed without error is obtained from the general model by fixing $\Lambda_y = I$, $\Lambda_x = I$, $\Theta_\delta = 0$ and $\Theta_\varepsilon = 0$.

The details of the formulations are provided in Hayduk (1989, chapter 4)

## Appendix 2

## Identification Problem in SEM:

We now illustrate the problem specifically in the context of SEM with a simple hypothetical example assuming that latent variables are perfectly correlated with the measurable variables, in other words in classical simultaneous equation model set-up. Suppose we consider a model with endogenous variables $y_1$, $y_2$ and an exogenous variable $x_1$. The model is,

$$y_1 = \gamma_{11} x_1 + \zeta_1, \qquad\qquad (A2.1)$$

$$y_2 = \beta_{21} y_1 + \zeta_2$$

where $\gamma_{11}$, $\beta_{21}$ are the regression coefficients and $\zeta_1$, $\zeta_2$ are the random error terms satisfying

Cov $(\zeta_1, x_1) = \text{Cov}(\zeta_1, \zeta_2) = \text{Cov}(\zeta_2, x_1) = 0$. Consistent with the notation introduced in (3.2) we have

$$B = \begin{bmatrix} 0 & 0 \\ \beta_{21} & 0 \end{bmatrix}, \quad \Gamma = \begin{bmatrix} \gamma_{11} \\ 0 \end{bmatrix}, \Psi = \begin{bmatrix} \psi_{11} & 0 \\ 0 & \psi_{22} \end{bmatrix}, \text{ and } \Phi = [\varphi_{11}]. \qquad (A2.2)$$

Consider now the equation obtained from (A2.1)

$$\mathrm{Cov}\begin{bmatrix} y_1 \\ y_2 \\ x_1 \end{bmatrix} = \mathrm{Cov}\begin{bmatrix} \gamma_{11}x_1 + \zeta_1 \\ \beta_{21}y_1 + \zeta_2 \\ x_1 \end{bmatrix},$$

which may equivalently be written as follows,

$$\begin{bmatrix} \mathrm{var}(y_1) & & \\ Cov(y_2, y_1) & \mathrm{var}(y_2) & \\ Cov(x_1, y_1) & Cov(x_1, y_2) & \mathrm{var}(x_1) \end{bmatrix} = \begin{bmatrix} \gamma_{11}^2\varphi_{11} + \psi_{11} & & \\ \beta_{21}(\gamma_{11}^2\varphi_{11} + \psi_{11}) & \beta_{21}^2(\gamma_{11}^2\varphi_{11} + \psi_{11}) + \psi_{22} & \\ \gamma_{11}\varphi_{11} & \beta_{21}\gamma_{11}\varphi_{11} & \varphi_{11} \end{bmatrix}. \qquad \text{(A2.3)}$$

Evidently the elements of the matrix on the left hand side of (A2.3) can be estimated from the sample variances and covariances of the observable variables. Also the sample values of the elements of the matrix on the left hand side represent the data constraints while the matrix on the right hand side represents the model constraints. Equating the elements of the matrices in (A2.3) component-wise we have six equations in five unknowns. Solving these we obtain the following estimates,

$\varphi_{11} = \mathrm{var}(x_1)$ , $\gamma_{11} = \mathrm{cov}(x_1, y_1)/ \mathrm{var}(x_1)$, $\beta_{21} = \mathrm{cov}(x_1, y_2)/ \gamma_{11}\varphi_{11}$

$\psi_{11} = \mathrm{var}(y_1) - \gamma_{11}^2\varphi_{11}$, or $[\mathrm{cov}(y_2, y_1)/\beta_{21}] - \gamma_{11}^2 \varphi_{11}$,

$\psi_{22} = \mathrm{var}(y_2) - \beta_{21}^2(\gamma_{11}^2 \varphi_{11} + \psi_{11})$. $\qquad\qquad$ (A2.4)

Here the model is said to be over-identified, since we have two sets of unique estimates of the parameters corresponding two choices of the estimate of $\psi_{11}$. In other words, here is a situation where we have more than one subset of constraints each leading to a unique set of estimates of the parameters. On the other hand, if the constraints lead to a single set of unique estimates, the model is called exactly identified. The above model becomes exactly identified if we assume that Cov $(\zeta_1, \zeta_2) = \psi_{12}$. If we augment the above model further by replacing $y_2 = \beta_{21} y_1 + \zeta_2$ with $y_2 = \beta_{21} y_1 + \gamma_{21} x_1 + \zeta_2$, the number of data constraints would then become less than the number of unknown parameters, and hence no unique set of estimates would be available. The model would then become under-identified or unidentifiable.

# Appendix 3

Estimation Methods of Structural Equation Parameters:

A measure of closeness between **S** and **Σ (θ)**, say, F(**S**, **Σ(θ)**) should satisfy the following natural conditions:

(i)       $F(\mathbf{S}, \Sigma(\boldsymbol{\theta}))$ is a scalar,

(ii)     $F(\mathbf{S}, \Sigma(\boldsymbol{\theta})) \geq 0$,

(iii)    $F(\mathbf{S}, \Sigma(\boldsymbol{\theta})) = 0$ iff $\Sigma(\boldsymbol{\theta}) = \mathbf{S}$,

(iv)    $F(\mathbf{S}, \Sigma(\boldsymbol{\theta}))$ is continuous in $\mathbf{S}, \Sigma(\boldsymbol{\theta})$.

Conditions (i)-(iii) are the properties that any measure of discrepancy should satisfy. The measure $F(\mathbf{S}, \Sigma(\boldsymbol{\theta}))$ is known as discrepancy function. A method of estimation is characterized by its choice of the discrepancy function. We give below the choices corresponding to the standard estimation methods,

Maximum Likelihood Method:

$F(\mathbf{S}, \Sigma(\boldsymbol{\theta})) = \log|\Sigma(\boldsymbol{\theta})| + \text{trace}[\mathbf{S}^{-1}\Sigma(\boldsymbol{\theta})] - \log|\mathbf{S}| - (p+q)$

Unweighted Least Squares Method:

$F(\mathbf{S}, \Sigma(\boldsymbol{\theta})) = 0.5 \, \text{trace}[\mathbf{S} - \Sigma(\boldsymbol{\theta})]^2$

Generalized Least Squares Method:

$F(\mathbf{S}, \Sigma(\boldsymbol{\theta})) = 0.5 \, \text{trace}[\{(\mathbf{S} - \Sigma(\boldsymbol{\theta}))\mathbf{W}^{-1}\}^2]$

The default choice of weight matrix W in almost all SEM software is S which reduces the discrepancy function to $F(\mathbf{S}, \Sigma(\boldsymbol{\theta})) = 0.5 \, \text{trace}[\{(\mathbf{I} - \Sigma(\boldsymbol{\theta}))\mathbf{S}^{-1}\}^2]$. The estimation method then uses an iterative procedure to minimize $F(\mathbf{S}, \Sigma(\boldsymbol{\theta}))$. If $\hat{\theta}$ minimizes $F(\mathbf{S}, \Sigma(\boldsymbol{\theta}))$ then it is taken as an estimate of $\boldsymbol{\theta}$.

## Glossary

**1. Latent variables:** Latent variables are hypothetical or unobserved variables. These are not directly observed but are rather captured using other observable variables.

**2. Construct (Concept):** A construct or equivalently a concept is an idea that unites phenomena like attitudes, behaviours, traits etc. under a single term. For instance the construct 'terrorism' provides the common element tying together diverse elements such as 'threat', 'use of violence' 'destruction of properties or lives of people' by individuals or groups for political purposes to shock or intimidate a target group wider than the immediate victims. The construct 'terrorism' acts as a summarizing device to replace a list of specific traits that an individual or a group may exhibit. Do constructs really exist? They are as real or as unreal as other ideas. They are created by people who believe that some phenomena have something in common. The measurement process begins with the definition of a construct.

**3. Observed (manifest) variables or Indicators:** Variables that can be directly measured or observed. It is the opposite of a latent variable, which cannot be directly observed. Manifest variables are used in measuring the latent variables. Models that connect the latent variable to the observed variables are called latent variable models. Manifest variables are considered either continuous or categorical.

**4. Theory:** A theory is an abstract set of ideas that links together constructs or concepts. For example we may desire to test a theory, "Democracy works as a deterrent to terrorism." Here the theory connects the two constructs, 'democracy' and 'terrorism'.

**5. & 6. Endogenous and Exogenous variables:** The terms endogenous and exogenous arise in the context of a model connecting several variables. A variable is called endogenous if it is explained within the model in which it appears. On the other hand a variable is called exogenous if it is determined by causes outside the model. For example the loyalty ($\xi_1$) of a customer to a soft drink brand is determined by trust ($\xi_2$) on the brand and the taste ($\xi_3$) of the customers. Trust is a variable determined by the model connecting loyalty to trust and taste, but taste is usually caused by factors outside the model. Thus trust and loyalty are endogenous while taste is exogenous.

**7. Path Diagrams:** It is a pictorial representation of a system of simultaneous equations. To understand a path diagram one needs to define the basic symbols used in such a diagram.
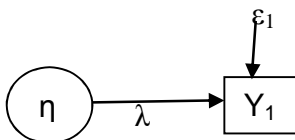
In the following we show it.

**Basic symbols used in Path Diagram**



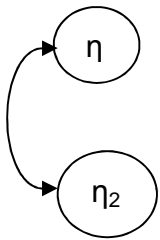Rectangular or square box represents an observed variable



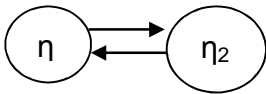Circle or ellipse represents a latent variable



Unenclosed variable represents an error term

Straight arrow signifies that variable at base of the arrow "causes" variable at head of arrow and $\lambda$ represents the regression of $Y_1$ on $\eta_1$

Curved two-headed arrow signifies assumed association between the two variables



Two single-headed arrows connecting two variables signifies Reciprocal causation

**8. & 9.Direct, indirect and total effects:** Path analysis classified into three types of effects: direct, indirect and total effects. The direct effect is the influence of one variable on another that is unmediated by any other variables in a path model. The indirect effect is provoked by at least one intervening variable. The sum of the direct and indirect effects is the total effect.

## 11. Exploratory Factor Analysis:

Let us consider the model

$X_i = \lambda_{i1} \xi_1 + \ldots + \lambda_{iq} \xi_q + \delta_i$ ,

where $X_i$'s ($i = 1,\ldots,p$) are the indicator variables and $\xi_i$'s are the latent variables representing the constructs or factors and $\delta_i$'s are uncorrelated random disturbance terms with variances $\sigma_i^2$'s. Defining $X = (X_1,\ldots,X_p)^T$, $\xi = (\xi_1,\ldots,\xi_q)^T$, $\lambda_i = (\lambda_{i1},\ldots,\lambda_{iq})^T, i = 1,\ldots, p$ $\Lambda = (\lambda_1,\ldots,\lambda_p)^T$ and Cov $(\xi) = I$ we have Cov($X$) $= \Lambda\Lambda^T + \text{diag}(\sigma_1^2, \ldots,\sigma_p^2)$. So unidimensionality is not achieved. However, by respecifying the model, sometimes unidimensionality could be achieved.

In an example, **Gerbing & Anderson (1988)** showed that exploratory factor analysis identifies two factors each substantially loading on five indicators. However, unidimensionality is achieved by removing two indicators from the model.

In exploratory factor analysis $\xi_1,\ldots,\xi_q$ are called the common factors and $\delta_i$'s are called the specific factors. If factor analysis works, we expect diag $(\sigma_1^2,\ldots,\sigma_p^2) \approx 0$, then Cov $(X) \approx \Lambda\Lambda^T$. Hence the covariance matrix of the $p$ indicator variables can be approximated by $q$ factors. Usually $q$ is much less than $p$. Now using spectral decomposition of Cov $(X)$, we have Cov $(X) = P$ diag$(\alpha_1,\ldots,\alpha_p)$ $P^T$ where $P = (P_1,\ldots,P_p)$, $P_1,\ldots,P_p$ are the eigen vectors and $\alpha_1,\ldots,\alpha_p$ are the eigen values of Cov $(X)$. Now if for $P^* = (P_1,\ldots,P_q)$ where $q<p$, Cov $(X) \approx P^* \text{diag}(\alpha_1,\ldots,\alpha_q)P^{*T}$ the factor loading matrix is given by $\Lambda = P^* diag(\alpha_1^{1/2},\ldots,\alpha_q^{1/2})$ .

## 10. Confirmatory Factor Analysis:

In exploratory factor analysis the researcher lets the data speak on the appropriate number of factors to extract along with the estimation of factor loadings. In confirmatory factor analysis on the other hand the researcher has theoretical reasons or past empirical evidence to believe they could predict the number of factors, and in extreme cases the actual values of those loadings. A confirmatory factor analysis presupposes the factor structure and thus specifies the measurement model. In our discussion, above we assume that the measurement model is pre-specified at the outset. The model is then estimated and finally we verify whether the data fit the model. In doing so often we need to revise the model due to its misspecification. In SEM this is often considered as confirmatory factor analysis. The analysis is carried out using standard software like LISREL and AMOS. However, some (Stewart, 2001) believe this is nothing but exploratory factor analysis camouflaged under the banner of confirmatory factor analysis. "Merely suggesting a structure and showing that data fit the suggested structure is not a genuine exercise in confirmatory factor analysis. An acceptable use of LISREL as a confirmatory tool requires at least three conditions:

1. A genuine, strong theory that posits a strong and unambiguous structure of relations among constructs and the variables that represent these constructs.

2. There must be a strong and unambiguous a priori structure that serves as the basis for the test of fit.

3. The fit of the data to the a priori structure must be better (by some acceptable criterion) than fit to the structure suggested by alternative theories; alternative structures that would be consistent with the theoretical foundation; intuitively obvious alternative structures; or structures that could be readily explained on methodological grounds, such as the presence of highly correlated error terms.

The others feel that variants of factor analysis should be placed in a continuum, with exploratory factor analysis on one end and confirmatory factor analysis in the strictest sense is at the other end.

## 12. Unidimensionality

For measurement of a construct we use more than one measures or proxies, considered to be alternative indicators of the same construct. A composite score corresponding to a respondent is generally calculated as an unweighted sum of the measures or proxies and is supposed to provide an estimate of the corresponding concept or construct. Computation of the composite score is meaningful if the measures are one-dimensional. "That a set of items forming an instrument all measure just one thing in common is most critical and basic assumption of measurement theory" (Hattie, 1985)

The mathematical definition of unidimensionality is based on the traditional common factor model in which a set of indicators $X_i$'s share only a single underlying factor $\xi$. Assuming linearity, the measurement model is given by

$$\mathbf{X_i = \lambda_i \, \xi_i + \delta_i,}$$

where $\boldsymbol{\lambda_i}$ is the factor loading and $\boldsymbol{\delta}_i$ is the random error.

## References:

1. Anderson, J. C, & Gerbing, D. W (1982). Some methods for respecifying measurement models to obtain unidimensional construct measurement. Journal of Marketing Research, 19, 453-460.

2. Anderson, J.C. (1987). Structural equation models in the social and behavioral sciences: Model building. Child Development,58, 49-64.

3. Anderson, J.C. (1987). An approach for confirmatory measurement and structural equation modeling of organized properties. Management Science,33, 525-541.

4. Anderson, J. C, & Gerbing, D. W., and Hunter, J.E. (1987). On the assessment of unidimensional measurement: Internal and external consistency, and overall consistency criteria. Journal of Marketing Research, 24, 432-437.

5. Anderson, J. C, & Gerbing, D. W (1988). Structural equation modeling in practice: A review and recommended two-step approach. Psychological Bulletin, 103, 411-423.

6. Bagozzi, R.P. (1981). Evaluating structural equation models with unobservable variables and measurement error: A comment., Journal of Marketing, 18, 375-381.

7. Bagozzi, R.P.(Edited) (1982). Special issue on causal modeling. Journal of Marketing, 19, 403-584.

8. Bagozzi, R.P.and Fornell, C. (1989). Consistency criteria and unidimensionality: An attempt at clarification. Advances in Consumer Research, 16, 321-325.

9. Brannick, M. T. (1995). Critical comments on applying covariance structure modeling. Journal of Organizational Behavior,16, 201-213.

10. Bauer, D.J. (2003).  The scaling of latent and observed variables in state-trait models. Measurement. Interdisciplinary Research and Perspectives, 1, 207-211.

11. Bentler, P.M. (1986). Structural modeling and Psychometrika: An historical perspective on growth and achievements. Psychometrika, 51, 35-51.

12. Bentler, P. M. and Tanaka, J. S. (1983). Problems with EM algorithms for ML factor analysis. Psychometrika, 48, 247-251.

13. Bentler, P.M. and Weeks, D.G. (1980). Linear Structural Equations with Latent Variables. Psychometrika, 45, 289 -308

14. Bentler, P. M. and Freeman, E. H. (1983). Test for Stability in Linear Structural Equation Systems. Psychometrika, 48, 143–145.

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

16. Bollen, K. A and Lennox, R. (1991). Conventional wisdom on measurement: A structural equation perspective. Psychological Bulletin, 110, 305-314.

17. Bullock, H. R. , Harlow, L. and Mulaik, S. A. (1994). Causation issues structural equation modeling research. Structural Equation Modeling, 1, 253-267.

18. Campbell, D. T. and Stanley, J. C. (1963). Experimental and quasi-experimental designs for research. Chicago: Rand McNally.

19. Cole, D.A. and Maxwell, S.E. (2003): Testing meditational models with longitudinal data: Questions and tips in the use of structural equation modeling, Journal of Abnormal Psychology, 112, 558-577.

20. Cliff, N. (1983). Some cautions concerning the application of causal modeling methods. Multivariate Behavioral Research, 18, 115-126.

21. Fornell, C. and Larker, D.F. (1981). Structural equation models with unobservable variables and measurement error: Algebra and statistics. Journal of Marketing Research, 18, 382-388.

22. Gerbing, D. W., & Anderson, J. C. (1984). On the meaning of within-factor correlated measurement errors. Journal of Consumer Research, 11, 572-580.

23. Gerbing, D.W. and Anderson, J.C. (1988). An updated paradigm for scale development incorporating unidimensionality and its assessment, Journal of Marketing Research, 186-192.

24. Gerbing, D.W. and Anderson, J.C. (1991). Predicting the performance of measures in a confirmatory factor analysis with a pretest assessment of their substantive validities, Journal of Applied Psychology, Vol. 76, No. 5, 732-740.

25. Grewal, R, Cote, J.A. and Baumgartner, H. (2004). Multicollinearity and measurement error in structural equation models: Implication for theory testing. Marketing Science, 23, 519-529.

26. Hattie, J.A (1985), Methodology review: Assessing unidimensionality of tests and items. Applied Psychological Measurement,9, 139–164.

27. Holden, R. R. and Jackson,D.N. (1979): Item subtlety and face validity in personality assessment. Journal of Consulting and Clinical Psychology, 47, 459-468.

28. Hoyle, R. H. and Smith, G. T. (1994). Formulating clinical research hypotheses as structural equation models: A conceptual overview. Journal of Consulting and Clinical Psychology, 62, 429-440.

29. James G. Anderson (1987). Structural Equation Models in the Social and Behavioral Sciences: Model Building. Society for Research in Child Development, 58, 49-64.

30. Jan-Benedict et al. (2000). On the use of structural equation models for marketing modeling. International Journal of Research in Marketing, 17, 195-202.

31. Kline, R.B. (1998). Principles and practice of structural equation modeling, The Guilford Press, New York.

32. Klein, L. (1950): Economic Fluctuations in the United States 1921—1941. New York: Wiley.

33. Krishnakumar, J. (2007). Going beyond functionings to capabilities: aneconometric model to explain and estimate capabilities, Journal of Human Development, 8, 39-63.

34. Krishnakumar, J. and A.L. Nagar. (2008). On exact statistical properties of multidimensional indices based on principal components, factor Analysis, MIMIC and structural equation models. Social Indicators Research, 86, 481-496.

35. Loevinger.J. (1957). Objective tests as instruments of psychological theory (monograph no. 9). Psychological Reports, 3, 635-694.

36. McArdle, J. J. (1980). Causal Modeling Applied to Psychonomic Systems Simulation. Behavior Research Methods and Instrumentation, 12, 193—209.

37. MacCallum, R. C., Wegener, D. T., Uchino, B.N. and Fabrigar, L.R. (1993). The problem of equivalent models in applications of covariance analysis. Psychological Bulletin, 114, 185-199.

38. Muthén, B. (1982). Some categorical response models with continuous latent variables. In K. G. Joreskog, & H. Wold (Eds.), Systems under indirect observation:Causality, Structure, Prediction (the conference volume). Amsterdam: North Holland.

39. Muthén, B. (1983). Latent variable structural equation modeling with categorical data. Journal of Econometrics, 22, 48-65.

40. Muthén, B. (1984). A general structural equation model with dichotomous, ordered categorical and continuous latent variable indicators. Psychometrika, 49, 115-132.

41. Muthén, B. (2002). Beyond SEM: General latent variable modeling.Behaviormetrika,29, 81-117.

42. Richard Netemeyer, Peter Bentler, Richard P Bagozzi, Robert Cudeck, Joseph Cote, Donald Lehmann, Roderick McDonald, Timothy Heath, Julie Irwin, Tim Ambler (2001). Structural equation modeling. Journal of Consumer Psychology,10, 83-100.

43. Armitage, P. and Colton, T. (2005). Encyclopedia of Biostatistics, John Wiley, 4363-4372.

44. Rabe-Hesketh, S., Skrondal, A. and Pickles, A.(2004). Generalized multilevel structural equation modeling. Psychometrika 69, 167-190.

45. Roderick P. McDonald, Moon-Ho Ringo Ho (2002). Principles and Practice in Reporting Structural Equation Analyses, Psychological Methods, Vol. 7, No. 1, 64–82.

46. Shash, R. and Goldstein, M.S. (2005). On the use of structural equation modeling in operations management research: Looking back and forward, Journal of Operations Management, 24, 148-169.

47. Skrondal, A. and S. Rabe-Hesketh (2004). Generalized Latent Variable Modeling: Multilevel, Longitudinal, and Structural Equation Models. Boca Raton: Chapman & Hall.

48. Sobel, M. E. (1993). Causal inferences in latent variable models. In Kenneth Bollen and John Long (Eds.). Testing structural equation models. 3-35, Newbury Park, CA: Sage.