# Structural Equation Modelling: A Powerful Antibiotic

## H. K. Dangi[1*], Ashmeet Kaur[2] and Juhi Jham[2]

[1]Associate Professor, Department of Commerce, Delhi School of Economics, University of Delhi, Delhi – 110021, India; hkd2students@gmail.com
[2]M.Phil Research Scholar, Department of Commerce, Delhi School of Economics, University of Delhi, Delhi – 110021, India; ashmeetkaur1810@gmail.com, juhi9294@gmail.com

## Abstract

This article is an attempt to scrutinize the applicability of the widely used statistical technique of Structural Equation Modelling (SEM). SEM is a comprehensive technique to test the model adequacy. SEM is considered as an important advancement in social science research as it combines measurement with substantive theories. It has been observed that many studies pay attention to statistical mechanisation of SEM rather than the assumptions on which it is based. The history of SEM can be traced to Regression Analysis, Path Analysis and Confirmatory Factor Analysis. SEM is popularly applied because of its use in estimating multiple dependence relationships. It is able to measure the unobserved variables, define the model representing the set of relationships and also corrects the measurement error. The technique is commonly applied in disciplines including sociology, psychology and other fields of behavioural science. The availability of various user-friendly software programmes like LISREL, AMOS, EQS, Mx, Mplus and PISTE is an added advantage. However, one should be careful while using SEM for causal inferences. In comparison to other common standard statistical techniques, SEM is based on several assumptions. The technique requires a priori knowledge of all the parameters to be estimated and a substantial amount of data pertaining to covariances, variances and path coefficients. It also requires relationships to be specified in the model. The model inherently assumes temporal precedence and is heavily dependent on researcher's judgements about exogeneity and directionality. Normality is yet another important assumption of SEM. The mismatch between data characteristics and assumptions imperils inference and accuracy. Like antibiotics are a boon to mankind yet one needs to judiciously use them. Similarly, SEM is a powerful technique however, researchers are suggested to apply cautiously.

**Keywords:** Confirmatory Factor Analysis, Latent Variables, Path Analysis, Regression Analysis, Structural Equation Modelling

**JEL classification:** C36, C38

## 1. Introduction

Structural Equation Modelling (SEM) is a statistical technique that subsumes and extends correlation, regression, factor analysis and path analysis. The technique focuses on the fit of the data to the theory. SEM is analysis of variance/covariance matrices of observed variables[4]. It yields an implied variance/covariance matrix which can be compared to observed matrix. It tests how well the constructed model fits the data as the model may be theoretically identified but empirically unidentified. It comprises testing two models namely measurement and structural model. Measurement model gives empirical evidences, while the structural model provides framework to support the hypotheses. SEM offers a major advantage of the latent variables being free of random error, leaving only a common variance. The technique involves building a model and starts with specification of a model. Model specification is probably the most critical as well as challenging part because it requires adequate knowledge and understanding of the theoretical models.

SEM is gaining popularity amongst researchers since it is a comprehensive technique that tests model adequateness. SEM has become an essential statistical tool and technique for academicians and practitioners. It is gaining popularity in its applications because of various reasons.

---

*Author for correspondence

Firstly, multiple observed variables can be studied simultaneously. It makes SEM efficient to deal with the modelling and testing of complex theories and phenomena. It also provides an edge to SEM over other basic statistical techniques using only limited number of variables. Secondly, greater relevance is given to the measurement error and the validity and reliability of the observed scores. This adds to the popularity of SEM techniques since the measurement error is taken explicitly into account while statistically analysing the data. Thirdly, several indices are reported for goodness and badness of model fit.

Over the past 30 years, SEM has matured in its ability of analysing and testing advanced theoretical models. In recent times, SEM has witnessed some major advancements such as the addition of new features of multilevel structural equation modelling, growth curve change modelling, generalized linear and mixed modelling, meta-analysis, and partial least squares. In those theoretical models wherein group differences are to be assessed, multiple-group SEM models can be used[11]. The development of multilevel SEM, which allows analysing data collected at more than one level (e.g. educational data) is another reason for its popularity.

All these advanced SEM models and techniques are an improvement over the basic statistical methods. A variety of SEM methods have been used in various fields such as business and sciences. But now, SEM techniques are widely used in the disciplines of biology, operation research, social, health and behavioural science. SEM is also evolving in the field of longitudinal investigations. All these improvements have enhanced the applications of SEM and are referred to as Second Generation SEM.

The various software programmes available for SEM are LISREL, AMOS, EQS, Mx, Mplus, PISTE. LISREL has been the most widely used programme since the 1970s. All these software programmes are user-friendly and allow the researchers to get the results conveniently.

The rest of the paper is organised as follows: Section 2 discusses the history behind the evolution of SEM from various statistical techniques. Section 3 reviews SEM as an important advancement in social science research and its applicability in various disciplines. Section 4 examines the rigid assumptions associated with the technique and how they limit its applicability. Concluding remarks on SEM being a powerful antibiotic are presented in Section 5.

# 2. Historical Perspective of SEM

It is important to trace the history of SEM. Basically the premise of SEM rests on four kinds of related models. They are; Regression Analysis, Path Analysis and Confirmatory Factor Analysis (CFA) and Structural Equation Models. Path Analysis and CFA are actually two special types of SEM. Regression model was developed by Francis Galton. It uses the criterion of Ordinary Least Squares and the coefficient of correlation to depict a relationship between two variables. Multivariate Regression models are used to depict the scores of the dependent observed variable (Y) through a set of independent observed variables (Xs) in a way that the sum of the squared residual error values is minimized. Structural Equation Models are different from regression models as they are relational and additive in nature.

Path Analysis is an extension of multiple regression[9]. It was first developed in the 1930s by Sewall Wright to be used in phylogenetic studies. It uses linear equation system to examine causal relationships between two or more variables. In the 1960s, Path Analysis was adopted in social sciences and since 1970s it is also being used in ecological sciences to a larger extent. In Path Analysis, we can have more than one dependent variable at a time. The variables can be both dependent and independent at the same time. There can be a chain of association wherein one variable can influence another variable, which, in turn, can influence the third variable. However, we substitute the terms independent and dependent variables with exogenous variables (which are not influenced by any other variables) and endogenous variables (which are influenced by other variables) respectively. A major limitation of Path Analysis is its inability to measure and analyse unobservable variables. Path Analysis is a subset of SEM. In SEM, we deal with unobservable variables, known as latent variables. These include the concepts which we encounter in everyday life such as anxiety, depression, quality of life and happiness. SEM extends Path Analysis and allows us to examine relations among the variables, both latent and measured.

The applications of factor analysis were developed by D. N. Lawley in the year 1940[8]. Over the years, factor analysis has been used in various disciplines to develop measurement instruments. It is being used to test theoretical models and constructs and assess whether the specified model fits the data or not. It is used when a

priori hypothesis is framed in regard to how the variables will cluster together on a factor[1]. It can be used to determine if the scale performs in the same manner while working with different population groups. The psychometric properties of varied versions of the scale can be compared using CFA. CFA is different from Exploratory Factor Analysis (EFA). In EFA, the software performs its statistical operations and produces the best combination of variables clustering together to form a factor, even if it is different from the combinations of variables that have been hypothesised. In those cases where the model does not fit the data, certain ideas and clues are available to guide the shuffling of variables so that the model fits the data in a better way. SEM is a combination of CFA and multiple in a broad sense.

## 3. Applications of SEM

SEM is a multivariate procedure that, as defined by Ullman[12], "Allows examination of a set of relationships between one or more independent variables, either continuous or discrete, and one or more dependent variables, either continuous or discrete." In SEM, the oval shape depicts latent variables while squares represent measured variables[10]. Each latent variable, also called construct or unobserved variable, is in fact a small CFA. Lines are used to indicate relationships between variables. Lines have either one arrow for depicting a hypothesised direct relationship between two variables or two arrows indicating a covariance between the two variables[13].

SEM is considered as an important advancement in social science research as it combines measurement with substantive theories. SEM analysis deals with testing of a model, testing a hypothesis about a model or modification of an existing model. The analysis makes it possible to simultaneously test all the relationships in case of complex and multidimensional constructs. SEM is commonly applied in disciplines including sociology, psychology and other behavioural science because of its capability to test relationships between latent and measured variables. Further, availability of many user-friendly software have increased the popularity of SEM amongst researchers.

To check model adequacy, goodness of fit test is used[5]. SEM is useful in understanding the relational data in multivariate systems and in examining the variances in the variables. It can be distinguished from other conventional methods of statistical analysis due to its distinct characteristics of using covariance as the basic statistic. Covariance statistic conveys more information than regression as in the latter, the differences between observed and expected individual cases are minimized. While, in SEM, the differences between observed and expected covariance matrices are minimized[14]. The analysis focuses on the fit of the data to the theoretical model. SEM allows us to distinguish between direct and indirect relationships among the variables by examining mediation and moderation.

SEM also indicates the group differences. Hence, it can be used to compare the results of separate models developed for different groups through multiple-group SEM models. The longitudinal data for measuring the change in the growth of variables over a period of time can be collected. It augments in providing new ideas to researchers in the field of SEM. The availability of user-friendly software programmes is an added advantage. Many of the software programmes are Windows-based and generate the programme syntax internally, and thus, are easier to apply.

SEM is popular in its applications because of the requisite use of multiple observed variables by researchers for better understanding. It is capable of dealing with the sophisticated theories that are statistically modelled and tested. All SEM have three main features[3]. First is the characteristic of estimating multiple dependence relationships. Second is the representation of unobserved concepts in the relationships and correction of the measurement errors. And third is the ability to define the model representing the set of relationships. Another contribution to its popularity is the measurement error taken explicitly into account while statistically analysing the data. This ensures a greater importance to the reliability and validity of the observed scores from the measurement instruments.

## 4. Limitations of SEM

SEM is a powerful technique for testing models but the modelling process, at times, is complicated. The technique requires a priori knowledge of all the parameters to be estimated. Significant amount of data pertaining to covariances, variances, path coefficients and the relationships is needed to be specified in the model. With the availability of user-friendly statistical softwares such as LISREL, AMOS and EQS, SEM is being used widely It is being excessively reported in social work journals without

adequately validating the assumptions on which the technique is fundamentally based. SEM inherently assumes temporal precedence i.e., presumed cause occurs before the presumed effect[6]. Temporal precedence must be checked through random assignment of cases to conditions in experimental studies and by measuring cause and effect relationships at different points over time in non-experimental studies. It is not possible to demonstrate temporal precedence if all variables are measured simultaneously, which is true in case of most of the studies, thus, rendering little justification for the inferences made.

It further assumes strong association or co-variation backed by both theories and results of empirical studies. This assumption significantly affects the inferences made as data may indicate spurious association. Association between two variables could be strong even if there is no causal relation as both variables might have a common third factor causing them. Prior knowledge of causal relations is assumed in interpreting path coefficients. In various fields of research like behavioural science, one barely knows the causal model, rather one hypothesizes the model. If the model fits the data one may conclude that model is consistent with the data but one cannot claim about the applicability of the model as it is not proven to be true.

Structural models are heavily dependent on the researcher's judgments about exogeneity and directionality. Exogeneity implies that the variable presumed to be exogenous must not affect the endogenous variable in any other way than prescribed, directly or indirectly. Such a variable must be uncorrelated with any other unmeasured cause of the endogenous variable. It is difficult to validate assumptions in the absence of robust empirical evidence.

Multivariate normality of the observed variable is yet another important assumption. Before building a model it is important to ensure that all the observations must be drawn from a continuous and multivariate normal population. The Maximum Likelihood (ML) technique of approximation, which assumes normality, is used to estimate the parameters[7]. This makes SEM a large sample technique. Thus, drawing conclusions from a small sample size makes them unreliable.

Use of non-continuous data in the model also leads to profusely biased results. Many studies have used this technique on dichotomous or ordinal data, which is an incorrect estimation method, as it may give inconsistent results. Estimation via ML technique further requires large sample size while sample sizes in social sciences research are mostly composed of less than a few hundred cases, thus, violating the basic assumption.

It has been observed that in many research studies, efforts made to model modifications for better goodness of fit statistic, are unnecessary. Significant modifications are made to improve a model fit. These include dropping indicators, allowing cross-loadings and including numerous correlated error terms, which not only challenge the reliability and validity of some studies, but also make it impossible to replicate the findings with new data.

Studies that make excessive modifications to a model also hamper meaningful interpretations as interpreting a model with numerous correlated error terms or cross-loadings or both, is inconclusive[2]. These modifications have improved the empirical fit of the model, however, theoretical consistency is compromised. Testing directionality in relationships is yet another challenge in the model. It is the researcher's hypotheses of causality that form the model. Recreation of the variance patterns, observed in nature, is not possible by using SEM since the working of the model is limited by the researcher's choice of variables and paths. This makes several models to fit the data equally well. Eliciting desired results is relatively easy indicating why the model is overly applied.

# 5. Conclusion

One should rather be careful while using SEM for causal inferences. Assumptions are critical in specification, analysis and interpretation. All statistical tests make certain assumptions about the data or model. It has been observed that inadequate attention is paid to assumptions. Mismatch between data characteristics and assumptions of a particular method used, imperils inference and accuracy of results. Conclusions that are extrapolated from a model based on a small sample size are unreliable. Omission of crucial variables is another major cause of poorly fit models.

SEM is being overly applied and hence, grossly misused. Unwarranted usage of SEM makes the validity of inferences questionable. It is being applied without an exhaustive knowledge of the variables and the model constructed, and without checking whether, or not, it is justified to use SEM. Further, the technique is deployed most of the times without validating the compliance of all the assumptions.

One may consider applying another suitable statistical test depending on the data properties. For non-normal

data, alternative methods of estimation like Ordinary Least Squares can be used with large sample size requirements. There should be no missing data in any variable. There are various methods dealing with such issues like Missing Completely at Random approach, Missing at Random approach and Imputation approach. One may use Partial Least Squares in place of AMOS for formative constructs. Thus, if assumptions are not met, suitable alternatives should be explored and employed.

Therefore, just like antibiotics, which are a big gift to mankind, but have been extensively misused resulting in drug resistance; SEM is also widely misused.

However, if used judiciously, it is a very powerful technique. Thus, there is a need for ensuring a more prudent use of SEM amongst researchers, warranting it is applied where apt and for the purpose it is meant for.

# 6. References

1. Anderson JC, Gerbing DW. Structural equation modeling in practice: A review and recommended two-step approach. APA PsycNET Direct. 1988; 103(3):411–23. https://doi.org/10.1037//0033-2909.103.3.411

2. Guo B, Perron BE, Gillespie DF. A systematic review of structural equation modelling in social work research. Br. J. Soc. Work 2009; 39(8):1556–74. https://doi.org/10.1093/bjsw/bcn101

3. Hair JF, Sarstedt M, Ringle CM, Mena JA. An assessment of the use of partial least squares structural equation modeling in marketing research. Journal of the Academy of Marketing Science. 2012; 40:414–33. https://doi.org/10.1007/s11747-011-0261-6

4. Hox JJ, Bechger TM. An introduction to structural equation modeling. Family Science Review. 1998; 11:354–73.

5. Kenny DA, McCoach DB. Effect of the number of variableson measures of fit in structural equation modeling. Structural Equation Modeling. 2003; 10(3):333–51. https://doi.org/10.1207/S15328007SEM1003_1

6. Kline RB. Assumptions in structural equation modeling. R. H. Hoyle, Ed. Handbook of structural equation modeling. The Guilford Press; 2012. p. 111–25.

7. Kumar S. Structural equation modeling basic assumptions and concepts: A novices guide. Asian Journal of Management Sciences. 2015; 03(07):25–8.

8. Lawley D. VI - The estimation of factor loadings by the method of maximum likelihood. Proceedings of the Royal Society of Edinburgh. 1940; 60(1):64–82. https://doi.org/10.1017/S037016460002006X

9. Streiner DL. Building a better model: An introduction to structural equation modelling. Can J Psychiatry. 2006; 51(5):317–24. PMid: 16986821. https://doi.org/10.1177/070674370605100507

10. Stoelting R. Structural Equation modeling/Path Analysis. 2002 Sept. http://userwww.sfsu.edu/efc/classes/biol710/path/SEMwebpage.htm

11. Tomarken AJ, Waller NG. Structural equation modeling: Strengths, limitations and misconceptions. 2005; 1:31–65. https://doi.org/10.1146/annurev.clinpsy.1.102803.144239

12. Ullman JB. Structural equation modeling. New York, NY: Harper Collins College Publishers; 1996. p. 709–819.

13. Ullman JB, Bentler PM. Structural equation modelling. Research Methods in Psychology. 2012; 2:663–83. https://doi.org/10.1002/9781118133880.hop202023

14. Valluzzi JL, Larson SL, Miller GE. Indications and limitations of structural equation modeling in complex surveys: Application in the Medical Expenditure Panel Survey (MEPS). Agency for Healthcare Research and Quality, Center for Financing, Access, and Cost Trends. 2003.