

Proceedings

Open Access

## Using simultaneous equation modeling for defining complex phenotypes

Terri M King\*

Address: Department of Internal Medicine, Division of Medical Genetics, The University of Texas – Houston Medical School, 6431 Fannin Street, Houston, Texas, USA

Email: Terri M King\* - Terri.M.King@uth.tmc.edu

\* Corresponding author

from Genetic Analysis Workshop 13: Analysis of Longitudinal Family Data for Complex Diseases and Related Risk Factors  
New Orleans Marriott Hotel, New Orleans, LA, USA, November 11–14, 2002

Published: 31 December 2003

*BMC Genetics* 2003, **4**(Suppl 1):S10

This article is available from: <http://www.biomedcentral.com/1471-2156/4/s1/S10>

### Abstract

**Background:** Interactions between multiple biological phenotypes are difficult to model. Simultaneous equation modelling (SEM), as used in econometric modelling, may prove an effective tool for this problem. Generalized linear models were used to derive the structural equations defining the interactions between cholesterol, glucose, triglycerides and high-density lipoprotein cholesterol (HDL-C). These structural equations were then applied, using SEM, to Cohort 2 data (replicates 1–100) to estimate the phenotypic structure underlying the simulation. The goal was to determine if this empiric method of deriving structural equations for use in SEM was able to recover the simulation model better than generalized linear models.

**Results:** First, the underlying structural equations were estimated using generalized linear model techniques, which found strong a relationship between glucose, triglycerides and HDL-C. Using these structural equations, I used SEM to evaluate these relationships jointly. I found that a combination of the empiric structural equations and the SEM method was better at recovering the underlying simulated relationship between biologic measures than generalized linear modelling.

**Conclusion:** The empiric SEM procedure presented here estimated different relationships between dependent variables than generalized linear modelling. The SEM procedure using empirically developed structural equations was able to recover the underlying simulation relationship partially and thus holds promise as a technique for complex phenotype analysis. Robust methods for determining the structural equations must be developed for application of SEM to population data.

### Background

To investigate complex relationships of interrelated phenotypes, I investigated whether simultaneous equation modelling (SEM) techniques can be used to detect this relationship in the absence of knowledge about the system. Simultaneous equation models describe two or more structural equations in which the dependent variable in one equation is a predictive variable in another. A classic

example from econometrics is the description of supply and demand within a population [1].

SEMs are attractive in longitudinal genetic studies because they have the ability to include 1) fixed data (e.g., genotypes), 2) variable data (e.g., cholesterol), and 3) stochastic data (e.g., cohort data) [2]. Using Problem Set 2 (complete data, replicates 1–100) without knowledge of

**Table 1: Results of Generalized Linear Models to Determine Structural Equations**

Dependent Variable	Independent Variable	Number of Replicates where the regression coefficient was:			Mean
		Significant	Not Significant	Normally Distributed?	
Cholesterol	<b>Age<sup>A</sup></b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.728</b>
	Cigarettes per day	2	98	Yes	-0.012
	Alcohol consumption	13	87	Yes	0.004
	Glucose	7	93	Yes	0.014
	HDL	10	90	Yes	-0.025
	Height	13	87	Yes	0.008
	<b>Systolic blood pressure</b>	<b>41</b>	<b>59</b>	<b>Yes</b>	<b>0.124</b>
	<b>Sex</b>	<b>43</b>	<b>57</b>	<b>Yes</b>	<b>3.916</b>
	Triglycerides	18	82	Yes	0.020
	Weight	11	89	Yes	0.007
Glucose	Age	14	86	Yes	-0.005
	Cigarettes per day	7	93	Yes	0.001
	Cholesterol	15	85	Yes	-0.016
	<b>Alcohol consumption</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>-0.081</b>
	<b>HDL</b>	<b>95</b>	<b>5</b>	<b>Yes</b>	<b>0.096</b>
	Height	5	95	Yes	-0.007
	Systolic blood pressure	22	78	Yes	0.020
	<b>Sex</b>	<b>90</b>	<b>10</b>	<b>Yes</b>	<b>2.075</b>
	<b>Triglycerides</b>	<b>100</b>	<b>0</b>	<b>No</b>	<b>0.081</b>
	<b>Weight</b>	<b>90</b>	<b>10</b>	<b>Yes</b>	<b>-0.002</b>
HDL-C	<b>Age</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.019</b>
	Cigarettes per day	10	90	Yes	-0.003
	<b>Cholesterol</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.160</b>
	<b>Alcohol consumption</b>	<b>100</b>	<b>0</b>	<b>No</b>	<b>0.249</b>
	<b>Glucose</b>	<b>95</b>	<b>5</b>	<b>Yes</b>	<b>0.111</b>
	Height	9	91	Yes	0.011
	Systolic blood pressure	10	90	Yes	-0.003
	<b>Sex</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>6.058</b>
	<b>Triglycerides</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>-0.194</b>
	<b>Weight</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.047</b>
Triglycerides	<b>Age</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>1.026</b>
	Cigarettes per day	18	82	Yes	0.012
	<b>Cholesterol</b>	<b>97</b>	<b>3</b>	<b>Yes</b>	<b>0.161</b>
	<b>Alcohol consumption</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.997</b>
	<b>Glucose</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.573</b>
	<b>HDL</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>-1.163</b>
	Height	8	92	Yes	0.018
	Systolic blood pressure	17	83	Yes	-0.002
	<b>Sex</b>	<b>100</b>	<b>0</b>	<b>No</b>	<b>-17.819</b>
	<b>Weight</b>	<b>100</b>	<b>0</b>	<b>Yes</b>	<b>0.268</b>

<sup>A</sup> Bolded variables were included in the structural equation.

the simulation structure, I examined a method to derive empirically the structural equations used in SEM to model the interrelationship phenotype at the first measurement point. The focus of this paper is to compare the ability to recover the simulation structure of traditional generalized linear models (GLM) to empirically derived structural equations used in SEM. This modelling technique would be useful in removing nongenetic components of variance prior to mapping efforts.

## Results

### Deriving structural equations

A summary of the GLM results is presented in Table 1. Using the results of the GLM analyses, the following structural equations were derived:

$$\text{Cholesterol} = c_1 + \alpha_1(\text{age}) + \alpha_2(\text{spb}) + \alpha_3(\text{sex}) + U_1 \quad (1)$$

$$\text{Glucose} = c_2 + \alpha_4(\text{HDL-C}) + \alpha_5(\text{sex}) + \alpha_6(\text{trig}) + \alpha_7(\text{wgt}) + U_2 \quad (2)$$

$$\text{HDL-C} = c_3 + \alpha_8(\text{age}) + \alpha_9(\text{cpd}) + \alpha_{11}(\text{gluc}) + \alpha_{12}(\text{sex}) + \alpha_{13}(\text{trig}) + \alpha_{14}(\text{wgt}) + U_3 \quad (3)$$

$$\text{Trig} = c_4 + \alpha_{15}(\text{age}) + \alpha_{16}(\text{cpd}) + \alpha_{17}(\text{drink}) + \alpha_{18}(\text{gluc}) + \alpha_{19}(\text{HDL-C}) + \alpha_{20}(\text{hgt}) + \alpha_{21}(\text{sex}) + \alpha_{22}(\text{wgt}) + U_4 \quad (4)$$

These results indicate the cholesterol is not a component of this system of structural models and thus was not included in the SEM analysis.

### Estimation of SEMs

Table 2 presents the results from the SEMs. The significant predictors of glucose in this system were: alcohol consumption, triglycerides, and weight. There were no significant predictors of high-density lipoprotein cholesterol (HDL-C). Finally, the significant predictors of triglycerides were: alcohol consumption, glucose, and weight. The direct generating variables in the simulation equations are denoted in italics in Table 2.

## Discussion

Fitting the GLMs consistently included more covariates than were used in the actual simulation equations. However, when the linear models were used to screen variable for structural equations and then SEM were used to determine the system, I was more successful in defining the underlying system.

The research presented here does not adequately address a number of key features that must be evaluated before endorsing this method. These include detection of non-linear and higher order relationships, the appropriate detection and adjustment of the correlation structure within the covariates, and estimation procedures in non-

replicate data. However, despite these elements being excluded from this research, I was encouraged by the ability of this method to provide a closer approximation to the simulation system than did GLMs.

## Conclusions

These results suggest that SEM can provide an alternative way to recover unknown relationships in complex phenotype data. The method presented here may result in a reduction in the model parameters that is overly conservative. Factors that must be evaluated in this relationship include the impact of the degree of correlation between the dependent and independent variables and the ability to detect a relationship with SEM.

This data structure seemed ideal to explore the usefulness of simultaneous equations for detailed deconstruction of complex phenotypes. This methodology, however, will need to overcome the challenges of defining robust structural models in the absence of knowledge of the underlying system. In this simulation study, I had the advantage of a large number of replicates, which, in real data, does not exist. I am currently investigating additional methods for determining the structural equations in undefined systems.

## Methods

### Data

Cohort 2 from replicates 1–100 was used in the complete data set without knowledge of the simulation conditions. The structural models were developed around four phenotypes at the first measurement time: cholesterol, glucose, HDL-C, and triglycerides, primarily because of literature focusing on the interrelationship of these agents [3,4]. Covariates included sex, age, height (hgt), systolic blood pressure (spb), cigarettes per day (cpd), alcohol consumption (drink), and weight (wt). Data were evaluated independent of familial structure. Covariates were tested and found to be normally distributed.

### Identification of the linear systems

The first step was to determine the structural equations that would be used in this analysis. To establish the structural equations, GLMs were fit in each of the 100 replicates to determine which of the covariates was significantly associated with each of the phenotypes. Using Proc GLM, within each replicate, the four phenotypes were analyzed with the following model structure.

$$\text{phenotype}_a = \text{intercept} + \text{phenotype}_b + \text{phenotype}_c + \text{phenotype}_d + \text{age} + \text{cpd} + \text{chol} + \text{drink} + \text{hgt} + \text{spb} + \text{sex} + \text{wgt} \quad (5)$$

Over the 100 replicates, the following information was collected on the regression coefficients: number of replicates in which the regression coefficient was significant ( $p$

**Table 2: Results of the Simultaneous Equations Model**

Dependent Variable	Independent Variable	Summary Statistics of the Regression Coefficients			
		Mean	Std Dev	Lower CI	Upper CI
Glucose	<b>Alcohol consumption<sup>A</sup></b>	<b>1.159</b>	<b>0.109</b>	<b>0.946</b>	<b>1.373</b>
	HDL	0.038	0.040	-0.040	0.116
	Sex	-0.196	0.110	-0.411	0.019
	<b>Triglycerides</b>	<b>-0.311</b>	<b>0.034</b>	<b>-0.378</b>	<b>-0.243</b>
	<i>Weight<sup>B</sup></i>	<i>0.276</i>	<i>0.009</i>	<i>0.258</i>	<i>0.294</i>
HDL	Age	-0.011	0.380	-0.755	0.734
	<i>Alcohol consumption</i>	<i>-5.436</i>	<i>186.730</i>	<i>-371.427</i>	<i>360.555</i>
	Glucose	2.408	176.659	-343.843	348.659
	Sex	2.470	21.777	-40.212	45.153
	Triglycerides	1.215	51.391	-99.512	101.942
	Weight	-0.384	50.526	-99.414	98.647
Tryglycerides	Age	-0.005	0.043	-0.089	0.078
	<b>Alcohol consumption</b>	<b>3.731</b>	<b>0.126</b>	<b>3.485</b>	<b>3.978</b>
	<b>Glucose</b>	<b>-3.176</b>	<b>0.433</b>	<b>-4.025</b>	<b>-2.327</b>
	<i>HDL</i>	<i>0.138</i>	<i>0.143</i>	<i>-0.142</i>	<i>0.418</i>
	Sex	-0.661	0.653	-1.940	0.619
	<b>Weight</b>	<b>0.877</b>	<b>0.114</b>	<b>0.563</b>	<b>1.285</b>

<sup>A</sup>Bolded variables were included in the structural equation. <sup>B</sup>Italicized variables were direct generators in the simulation model.

< 0.05), the average of regression coefficient, and whether the distribution of the regression coefficient was normally distributed. To establish the structural equations, covariates were selected that had regression coefficients that were significant in more than 25% of the replicates. It is important to note that I was not interested in the value of the regression coefficient *per se*, but rather if that regression coefficient was significant in a percentage of GLM models.

### Estimation of equations

Using equations (1–4) above, the associated parameters ( $\alpha_4 - \alpha_{22}$ ) were estimated using Proc Syslin within SAS [5] for each replicate. For this analysis, the parameters were estimated using two-stage least-squares techniques, which allow for these violations. In these techniques, the models are restructured with temporary dependent variables that are not in violation of the recursivity assumption. Then the models are estimated using ordinary least-square methods. These results are presented in Table 2.

### References

1. Goldberger AS: **Introductory Econometrics**. Cambridge, MA, Harvard University Press 1998.
2. Wooldridge JM: **Econometric Analysis of Cross Section and Panel Data**. Cambridge, MA, The MIT Press 2002.
3. Bosselo O, Zamboni M: **Visceral obesity and metabolic syndrome**. *Obes Rev* 2000, **1**:47-56.

4. Knopp RH: **Risk factors for coronary artery disease in women**. *Am J Cardiol* 2002, **89**(12 suppl):28E-34E. discussion 34E-35E
5. The SAS Institute Inc.: **Statistical Analysis Software v8.1**. Cary, NC, SAS Institute, Inc. .

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:  
[http://www.biomedcentral.com/info/publishing\\_adv.asp](http://www.biomedcentral.com/info/publishing_adv.asp)

