

LATENT VARIABLE MODELS

Session 1 – Introduction

1

Course prerequisites

2

- Knowledge of statistics
 - ▣ Statistical testing (e.g., chi-square & normal distributions)
 - ▣ Regression (GLMs)
 - ▣ Var, cov, cor, mean
- Knowledge of psychometrics
 - ▣ Validity
 - ▣ PCA, EFA, CFA
 - ▣ Reliability
 - ▣ IRT
- Programming in R

2

Course materials

3

Book(s):

- Beaujean, A. A. (2014). *Latent variable modeling using R: A step-by-step guide*.
 - ▣ Good as a starting guide, not an authoritative standard
- Kaplan, D. (2009). *Structural Equation Modeling: Foundations and Extensions*.
 - ▣ Authoritative standard. But more technical and not focused on specific software, so less practical.

Brightspace materials:

- Lecture slides
- Markdown files for examples and exercises

3

Book examples

4

- You are strongly advised to copy and run R code from examples in Beaujean book and from Brightspace:
 - ▣ They give you a step-by-step guide on how to perform analyses
 - ▣ They give you a starter for making the exercises
 - ▣ If you make a mistake, you will get an error or warning message, from which you learn A LOT! (But only you read and try to decipher! (red = good!))

4

Structural Equation Modeling

5

SEM: the modeling of structural equations




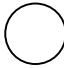

- ▣ **Modeling:** we construct models (hypotheses, theories) of reality. The models (theory, hypothesis) can be statistically tested. That is: rejected by the data (or not), but never proven 'true' or 'right'.
 - In fact, all models are wrong, but some are useful.
- ▣ **Structural:** the model is used to explain the interrelations between (that is, the structure of) observed variables
- ▣ **Equations:** the interrelations between variables in the model are described using mathematical formulae (equations)

5

Structural Equation Modeling

6

SEMS are graphically represented using these building blocks:

- | | |
|---|--|
|  | Observed (manifest) variable |
|  | Directional relationship (regression relationship) |
|  | Non-directional relationship
(correlation/(co)variance) |
|  | Latent variable |
|  | Constant term (i.e., not a variable, e.g., intercept) |

6

Structural Equation Modeling

7

- The arrows in SEM denote regression relationships
- All generalized linear models (GLMs) can be formulated as SEM models:
 - t-test
 - ANOVA
 - Multiple linear regression
 - Multiple logistic regression
 -
- Also, SEM can be used to models for multilevel or longitudinal data (i.e., GLMMs)

7

Example dataset

8

Variables in the model:

grade

ethnicity

homework

SES

Prev_ach

- GPA in 10th grade
- Ethnicity
- Homework (8th grade)
- Socio-economic status
- Previous achievement (8th grade)

□ Sample covariance matrix **S**:

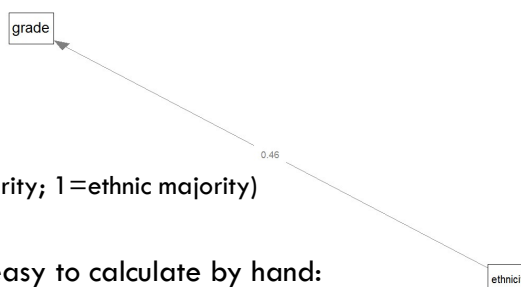
	grade	homwrk	prv_ch	ethnct	SES
grade	2.185				
homework	0.335	0.649			
prev_ach	6.429	2.067	79.092		
ethnicity	0.081	0.028	1.201	0.175	
SES	0.338	0.176	3.541	0.106	0.690

8

Model: Univariate regression

9

- Dependent:
 - ▣ GPA in 10th grade
- Independent:
 - ▣ ethnicity (0=ethnic minority; 1=ethnic majority)
- Regression coefficient easy to calculate by hand:
 - ▣ $\hat{b}_{xy} = \frac{cov_{x,y}}{var_x} = \frac{0.0814}{0.1752} = 0.4646$
 - ▣ standardized $\hat{b}_{xy} = \hat{\rho}_{xy} = \frac{cov_{x,y}}{s_x s_y} = s_x \frac{\hat{b}_{xy}}{s_y} = 0.132$
- Measure of fit or (strength of) association: $\hat{\rho}_{xy}^2$

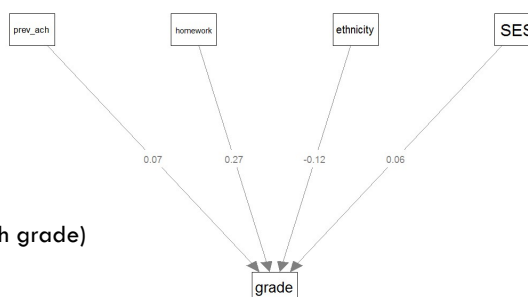


9

Model: Multiple regression

10

- Dependent:
 - ▣ GPA in 10th grade
- Independent:
 - ▣ Ethnicity
 - ▣ Homework (8th grade)
 - ▣ Previous achievement (8th grade)
 - ▣ Socio-economic status
- Regression estimates are now a vector of partial regression coefficients, need matrix algebra to compute: $\hat{\beta} = (X^T X)^{-1} X^T y$
- Measure of fit: multiple correlation ($R=.512$), or variance explained ($R^2=.262$)
- Measure of (strength of) association: \hat{b}_{xy} or standardized $\hat{b}_{xy} = s_x \frac{\hat{b}_{xy}}{s_y}$ (where \hat{b}_{xy} is now a partial regression coefficient)

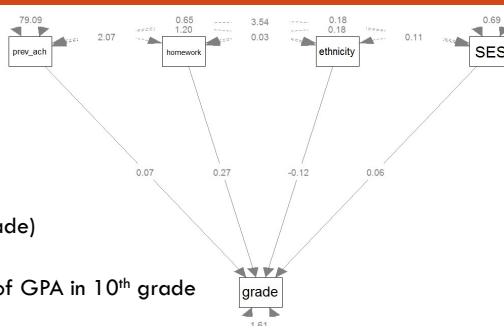


10

Model: SEM

11

- Endogenous variables:
 - GPA in 10th grade
- Exogenous variables:
 - Ethnicity
 - Homework (8th grade)
 - Previous achievement (8th grade)
 - Socio-economic status
 - (disturbance/error/residual) of GPA in 10th grade
- Regression estimates are still a vector of partial regression coefficients, need matrix algebra and optimization to compute
- Measure of (strength of) associations: Partial regression coefficients
- Overall model fit: How well are the observed variables' (co)variances reproduced by the model?
 - Quantified by a χ^2 value and model fit indices

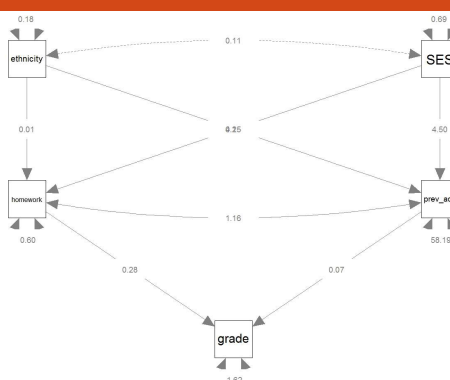


11

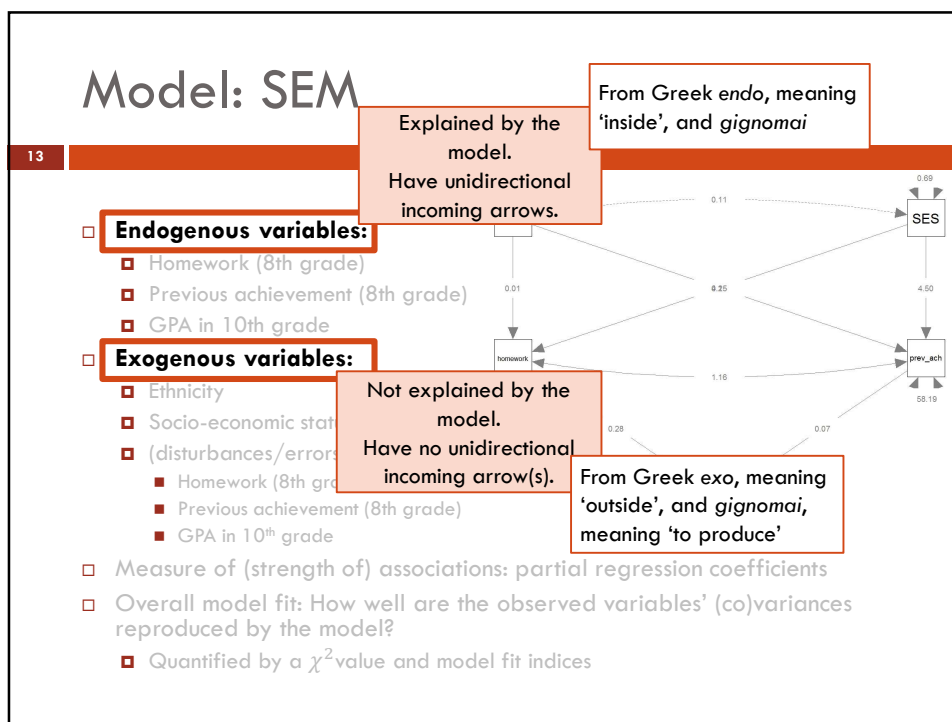
Model: SEM

12

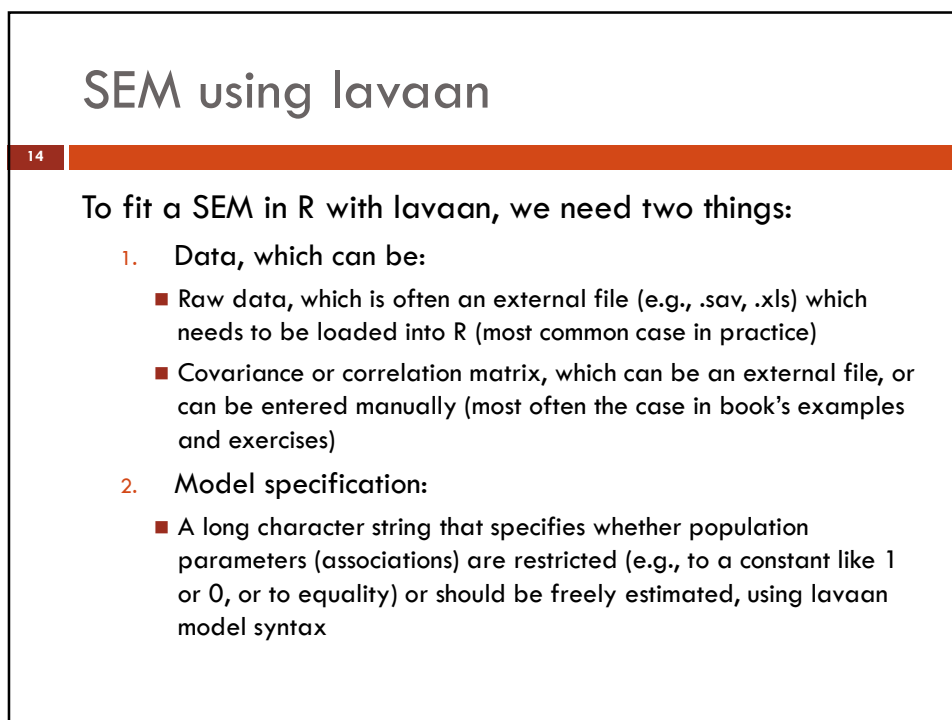
- Endogenous variables:
 - Homework (8th grade)
 - Previous achievement (8th grade)
 - GPA in 10th grade
- Exogenous variables:
 - Ethnicity
 - Socio-economic status
 - (disturbances/errors/residuals of
 - Homework (8th grade)
 - Previous achievement (8th grade)
 - GPA in 10th grade
- Measure of (strength of) associations: partial regression coefficients
- Overall model fit: How well are the observed variables' (co)variances reproduced by the model?
 - Quantified by a χ^2 value and model fit indices



12



13



14

Lavaan model syntax

15

Syntax	Command	Example
~	Regress onto	Regress B onto A: $B \sim A$
~~	(Co)variance	Variance of A: $A \sim\sim A$ Covariance of A and B: $A \sim\sim B$
~1	Constant/mean/intercept	Regress B onto A, and include the intercept in the model: $B \sim 1 + A$ or $B \sim A$ $B \sim 1$
=~	Define reflective latent variable	Define Factor 1 by A-D: $F1 =\sim A+B+C+D$
<~	Define formative latent variable	Define Factor 1 by A-D: $F1 <\sim 1*A+B+C+D$
:=	Define non-model parameter	Define parameter u2 to be twice the square of u: $u2 := 2*(u^2)$
*	Label parameters (the label has to be pre-multiplied)	Label the regression of Z onto X as b: $Z \sim b*X$
	Define the number of thresholds (for categorical endogenous variables)	Variable u has three thresholds: $u t1 + t2 + t3$

15

grade

ethnicity

homework

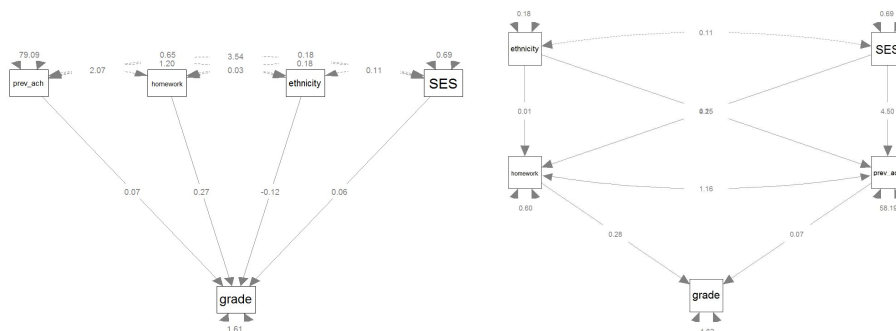
SES

Prev_ach

16

Syntax	Command	Example
~	Regress onto	Regress B onto A: $B \sim A$
~~	(Co)variance	Variance of A: $A \sim\sim A$ Covariance of A and B: $A \sim\sim B$

Q: How do we specify these models in lavaan syntax?



16

Computation time!

17

Example 2.4.1

- get PDF from Github



Make Exercise 2.1:

- Get Exercises_week_1.pdf from Github (adapted version of the exercises in the Beaujean book)

17

Structural Equation Modeling

18

- Fitted model is used to **explain the structure** of, or the interrelations between observed variables
- That is, to explain covariances between observed variables:

$$\text{cov}_{xy} = \left(\frac{1}{N-1} \right) \sum_i (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\text{cov}_{xy} = r_{xy} SD_x SD_y$$

- Note: means, and skewness & kurtosis can be also be involved in SEM (discussed later in course)

18

Structural Equation Modeling

19

- With SEM, we obtain a fitted model that minimizes the difference between
 - ▣ sample matrix of observed covariances \mathbf{S} and
 - ▣ population matrix of model-implied covariances $\hat{\Sigma}$
 - In addition, we try to keep the model parsimonious through applying restrictions (i.e., specifying the model) so that not all possible paths are estimated
- These covariance matrices contain all (co)variances of the observed variables in the model. Note that:
 - ▣ Covariance matrices are always symmetric, because $\text{cov}(x,y) = \text{cov}(y,x)$
 - ▣ Covariance matrices have the variance of the observed variables on the diagonal. I.e., $\text{cov}(x,x) = \text{var}(x)$

19

Model-implied (co)variances

20

- Variables in the model:



- Observed covariance matrix \mathbf{S} :

```

      grade  homwrk prv_ch ethnct SES
grade      2.185
homework   0.335  0.649
prev_ach   6.429  2.067 79.092
ethnicity  0.081  0.028 1.201 0.175
SES        0.338  0.176 3.541 0.106 0.690
  
```

- Once the model is estimated, the model-implied covariance matrix $\hat{\Sigma}$ can be calculated using path analysis, or equivalently, matrix algebra

20

Model-implied (co)variances

21

Path analysis:

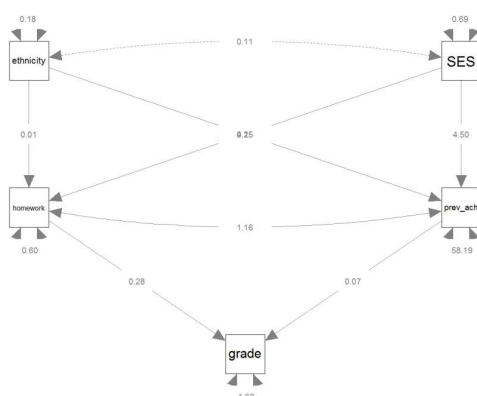
- Model-implied covariance between variables X and Y can be computed as follows:
 - Find all paths leading from X to Y
 - Multiply all parameter values along a given path from X to Y, but:
 - No loops: may not go through same variable more than once
 - May switch forward/backward direction only once within a path
 - May go through double-headed arrow only once within a path
 - Summing all values thus obtained
- Variances of variables are calculated as follows:
 - For exogenous variables, model-implied variances are equal to sample variances, so are given (not computed)
 - For endogenous variables, variances are computed like covariances (rules above)

21

Model-implied (co)variances

22

Model:



Parameter estimates:

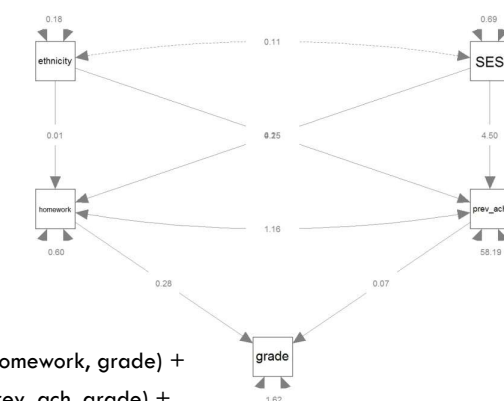
lhs	op	rhs	est
grade	~	prev_ach	0.074
grade	~	homework	0.281
homework	~	ethnicity	0.007
homework	~	SES	0.254
prev_ach	~	ethnicity	4.147
prev_ach	~	SES	4.496
homework	~	prev_ach	1.158
grade	~	grade	1.616
homework	~	homework	0.604
prev_ach	~	prev_ach	58.190
ethnicity	~	ethnicity	0.175
ethnicity	~	SES	0.106
SES	~	SES	0.690

model-implied cov(SES, grade)?

22

Model-implied (co)variances

23



model-implied $\text{cov}(\text{SES}, \text{grade}) =$

$$\begin{aligned} & \text{var}(\text{SES}) * b(\text{SES}, \text{homework}) * b(\text{homework}, \text{grade}) + \\ & \text{var}(\text{SES}) * b(\text{SES}, \text{prev_ach}) * b(\text{prev_ach}, \text{grade}) + \\ & \text{cov}(\text{SES}, \text{ethnicity}) * b(\text{ethnicity}, \text{homework}) * b(\text{homework}, \text{grade}) + \\ & \text{cov}(\text{SES}, \text{ethnicity}) * b(\text{ethnicity}, \text{prev_ach}) * b(\text{prev_ach}, \text{grade}) \end{aligned}$$

23

Model-implied (co)variances

24

model-implied $\text{cov}(\text{SES}, \text{grade}) =$

$$\begin{aligned} & \text{var}(\text{SES}) * b(\text{SES}, \text{homework}) * b(\text{homework}, \text{grade}) + \\ & \text{var}(\text{SES}) * b(\text{SES}, \text{prev_ach}) * b(\text{prev_ach}, \text{grade}) + \\ & \text{cov}(\text{SES}, \text{ethnicity}) * b(\text{ethnicity}, \text{homework}) * b(\text{homework}, \text{grade}) + \\ & \text{cov}(\text{SES}, \text{ethnicity}) * b(\text{ethnicity}, \text{prev_ach}) * b(\text{prev_ach}, \text{grade}) = \end{aligned}$$

$$\begin{aligned} & .690 * .254 * .281 + \\ & .690 * 4.496 * .074 + \\ & .106 * .007 * .281 + \\ & .106 * 4.147 * .074 = \\ & 0.3115514 \end{aligned}$$

Note that Beaujean's examples in section 2.1.3 seem more simple, because he uses the standardized solution. Then all variances of exogenous variables equal 1 and can be omitted, which simplifies calculations a lot.

lhs	op	rhs	est
grade	~	prev_ach	0.074
grade	~	homework	0.281
homework	~	ethnicity	0.007
homework	~	SES	0.254
prev_ach	~	ethnicity	4.147
prev_ach	~	SES	4.496
homework	~	prev_ach	1.158
grade	~	grade	1.616
homework	~	homework	0.604
prev_ach	~	prev_ach	58.190
ethnicity	~	ethnicity	0.175
ethnicity	~	SES	0.106
SES	~	SES	0.690

24

Model-implied (co)variances

25

- A SEM is a system of linear equations, which we can represent by matrices
 - ▣ Although non-linear SEM also exists, but outside the scope of this course
- The tracing rules represent matrix algebra but more tedious/confusing/error prone
- Beaujean's book hardly involves formulas, and no matrix notation. To get a good understanding of SEM, you need to know about underlying matrices and vectors

25

Model-implied (co)variances

26

- In lavaan, the (co)variance structure of a fitted model is given by four parameter matrices
- Matrix algebra gives us the model-implied covariance matrix:

$$\hat{\Sigma} = \Lambda(\mathbf{I} - \beta)^{-1} \Psi[(\mathbf{I} - \beta)^{-1}]^T \Lambda^T + \Theta$$

- Today, our models assume no measurement error, so Λ is an identity matrix and Θ all zeros. Thus, the above formula simplifies to:

$$\hat{\Sigma} = (\mathbf{I} - \beta)^{-1} \Psi[(\mathbf{I} - \beta)^{-1}]^T$$

26

Model-implied (co)variances

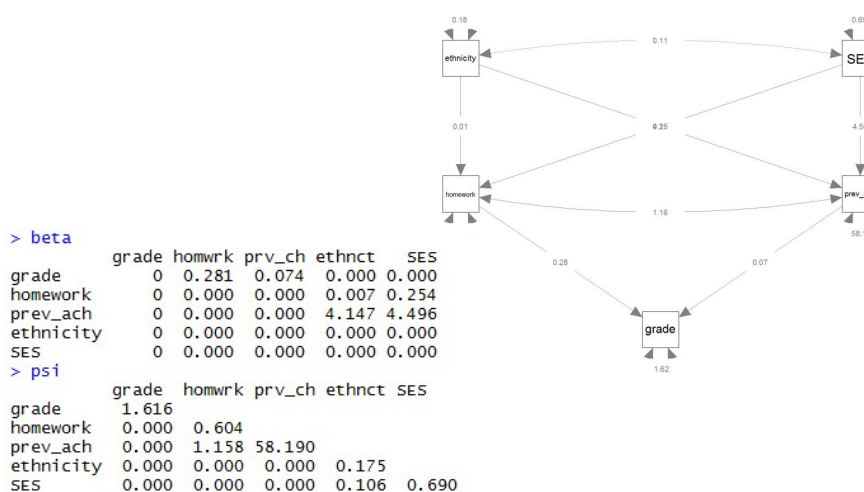
27

- Let p be the number of observed variables in the model
- If we have observed variables only:
 - β is a $p \times p$ matrix of regression coefficients, relating predictor to criterion variables
 - 'Contains' single-headed (directed) arrows, therefore non-symmetric
 - The columns reflect the variables as predictors, the rows reflect the variables as responses
 - ψ is a $p \times p$ matrix of (co)variances not explained by the regression equations
 - 'Contains' double headed (undirected) arrows, therefore symmetric
- ψ and β describe the **structural** model
- Often, SEM models also involve a **measurement** model (described by Λ and Θ , which will be introduced next session)

27

Model-implied (co)variances

28



28

Structural and measurement model

29

Two main components of SEMs:

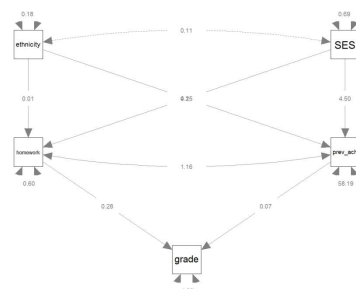
- ▣ the ***structural model*** contains *causal* regression relationships between endogenous and exogenous variables
 - path models (without measurement errors) can be viewed as SEMs that contain only the structural model
- ▣ the ***measurement model*** contains the associations between latent variables and their indicators
 - confirmatory factor analysis models contain only the measurement part

29

Model-implied (co)variances

30

Model:



Parameter estimates:

lhs	op	rhs	est
grade	~	prev_ach	0.074
grade	~	homework	0.281
homework	~	ethnicity	0.007
homework	~	SES	0.254
prev_ach	~	ethnicity	4.147
prev_ach	~	SES	4.496
homework	~	prev_ach	1.158
grade	~	grade	1.616
homework	~	homework	0.604
prev_ach	~	prev_ach	58.190
ethnicity	~	ethnicity	0.175
ethnicity	~	SES	0.106
SES	~	SES	0.690

Model-implied covariance matrix $\hat{\Sigma}$:

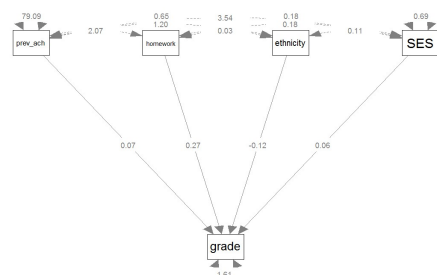
	grade	homwrk	prev_ch	ethnct	SES
grade	2.185				
homework	0.335	0.649			
prev_ach	6.429	2.067	79.092		
ethnicity	0.097	0.028	1.201	0.175	
SES	0.311	0.176	3.541	0.106	0.690

30

Model-implied (co)variances

31

Model:



Parameter estimates:

lhs	op	rhs	est
grade	~	prev_ach	0.073
grade	~	homework	0.271
grade	~	ethnicity	-0.119
grade	~	SES	0.063
grade	~~	grade	1.612
prev_ach	~~	prev_ach	79.092
prev_ach	~~	homework	2.067
prev_ach	~~	ethnicity	1.201
prev_ach	~~	SES	3.541
homework	~~	homework	0.649
homework	~~	ethnicity	0.028
homework	~~	SES	0.176
ethnicity	~~	ethnicity	0.175
ethnicity	~~	SES	0.106
SES	~~	SES	0.690

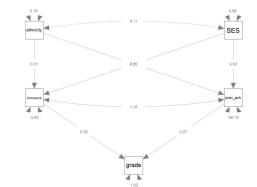
Model-implied covariance matrix $\hat{\Sigma}$:

	grade	prv_ch	homwrk	ethnct	SES
grade	2.185				
prev_ach	6.429	79.092			
homework	0.335	2.067	0.649		
ethnicity	0.081	1.201	0.028	0.175	
SES	0.338	3.541	0.176	0.106	0.690

31

Model-implied (co)variances

32

 $\hat{\Sigma} =$

	grade	homwrk	prv_ch	ethnct	SES
grade	2.185				
homework	0.335	0.649			
prev_ach	6.429	2.067	79.092		
ethnicity	0.097	0.028	1.201	0.175	
SES	0.311	0.176	3.541	0.106	0.690

 $\hat{\Sigma} =$

	grade	prv_ch	homwrk	ethnct	SES
grade	2.185				
prev_ach	6.429	79.092			
homework	0.335	2.067	0.649		
ethnicity	0.081	1.201	0.028	0.175	
SES	0.338	3.541	0.176	0.106	0.690

 $S =$

	grade	homwrk	prv_ch	ethnct	SES
grade	2.185				
homework	0.335	0.649			
prev_ach	6.429	2.067	79.092		
ethnicity	0.081	0.028	1.201	0.175	
SES	0.338	0.176	3.541	0.106	0.690

Which model fits data best (i.e., approximates sample covariances best)?

Which is most parsimonious (i.e., estimates lowest number of population parameters)?

32

Variances of exogenous variables often not explicitly depicted

33

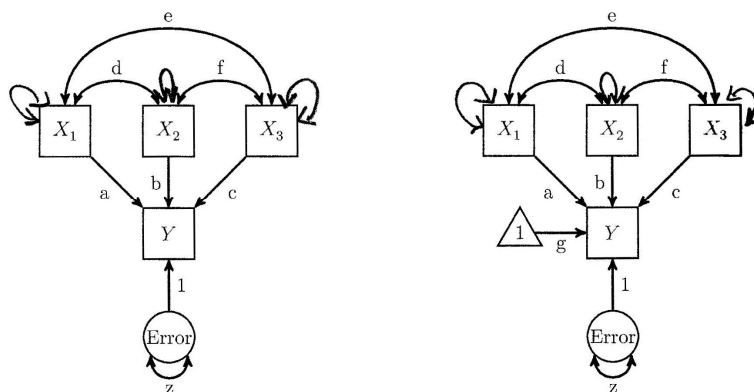
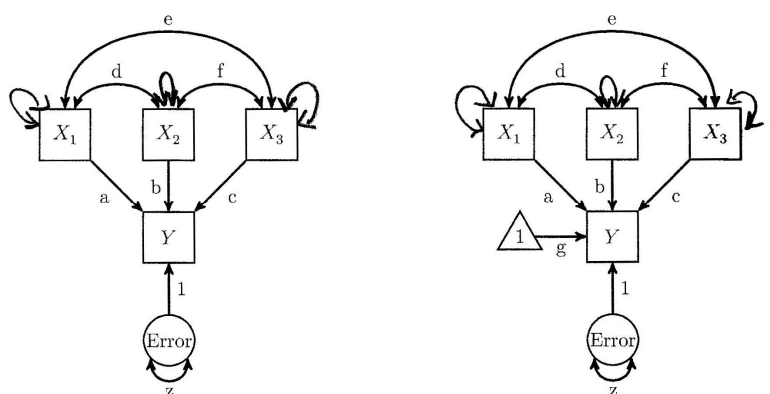


Figure 2.2 Path model of a multiple regression with three predictor (exogenous) variables.

33

Mean structure often omitted

34



(co)variance structure only
all means omitted (i.e., assumed zero)
 $Y = aX_1 + bX_2 + cX_3 + error$

(co)variance and mean structure
means freely estimated
 $Y = g + aX_1 + bX_2 + cX_3 + error$

Figure 2.2 Path model of a multiple regression with three predictor (exogenous) variables.

34

Error terms

35

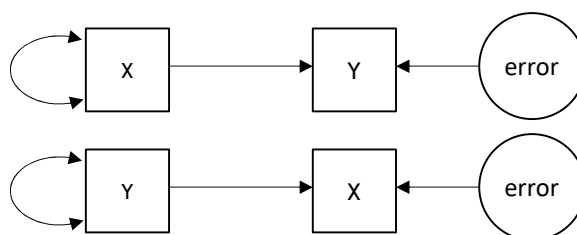
- Errors are also latent variables: they are hypothetical, not directly observed
- Error is defined as the difference between observed (sample) variance and variance explained by other variables in the model
 - ▣ Therefore, a variable that has an error/disturbance term is an endogenous variable
 - ▣ Errors/disturbance terms are always exogenous (have no incoming directional arrows)

35

Causation

36

- Causation is a function of the research design, and cannot be determined statistically



- Both models above will fit the observed data equally well, it is up to the researcher to decide on the direction of the arrows!
 - ▣ In the SEM model, it is merely a matter of scaling:

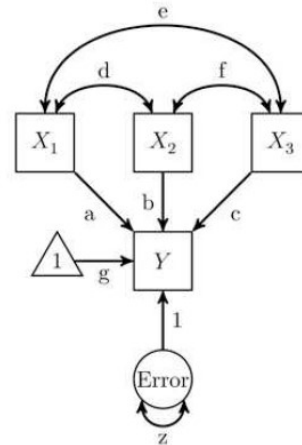
$$b_x = \frac{cov_{xy}}{var_x} \text{ and } b_y = \frac{cov_{xy}}{var_y}$$

36

Path & partial regression coefficients

37

- Path coefficients (a , b , c , g and 1) are partial regression coefficients
- That is, the expected increase in the response variable, when the predictor variable increases by 1, controlling for (= keeping constant) all the other predictor variables
- Note that the intercept is always 1, so cannot increase or decrease



37

Standardized coefficients

38

- Parameter estimates (path coefficients) can be standardized and unstandardized
 - Unstandardized: Interpret like regression coefficients
 - Expected increase in Y if X increases by 1
 - Standardized: Interpret like correlation coefficients
 - Expected increase in SDs of Y if X increases by 1 SD
 - 0: no linear association; -1: perfect negative association; 1: perfect positive association
 - squared standardized coefficient = prop. of variance in Y explained by X (vice versa)

38

Lavaan model syntax

39

Syntax	Command	Example
~	Regress onto	Regress B onto A: $B \sim A$
~~	(Co)variance	Variance of A: $A \sim\sim A$ Covariance of A and B: $A \sim\sim B$
~1	Constant/mean/intercept	Regress B onto A, and include the intercept in the model: $B \sim 1 + A$ or $B \sim A$ $B \sim 1$
=~	Define reflective latent variable	Define Factor 1 by A-D: $F1 =\sim A+B+C+D$
<~	Define formative latent variable	Define Factor 1 by A-D: $F1 <\sim 1*A+B+C+D$
:=	Define non-model parameter	Define parameter u2 to be twice the square of u: $u2 := 2*(u^2)$
*	Label parameters (the label has to be pre-multiplied)	Label the regression of Z onto X as b: $Z \sim b*X$
	Define the number of thresholds (for categorical endogenous variables)	Variable u has three thresholds: $u t1 + t2 + t3$

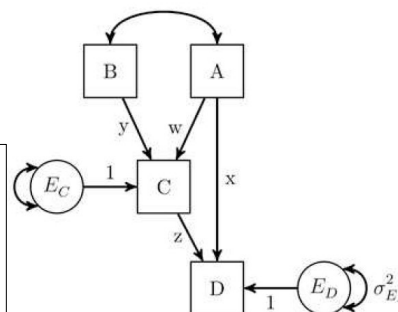
39

Lavaan syntax exercise

40

- 1) How do we write the model below in lavaan syntax?
- 2) How can we label and refer to the indirect effect from A on D via C in lavaan syntax?
- 3) What do the beta and psi matrices for this model look like?

Note that Beaujean often labels paths in lavaan syntax, but that is not required - I never do it, unless there are indirect effects that I want to explicitly define in the model. It does not make a difference for the estimated parameters and model fit.



40

Homework

41

- Exercises 2.2 and 2.3 (see PDF on Brightspace)
- See Example-2.4.1.pdf on Brightspace for instructions on extracting beta and psi matrices

41