

SEM 2: Structural Equation Modeling

Week 1 - Causal modeling and SEM

Sacha Epskamp

03-05-2018

Course Overview

- Tuesdays: Lecture
- Thursdays: Unstructured practicals
- Three assignments
 - First two 20% of final grade, last 10% of final grade
- Final project
 - Presentations and report, 50% of final grade

Schedule

“Week 1” – Introduction to Structural Equation Modeling

- Thursday May 3 – Lecture + practical
- Tuesday May 8 – Lecture + practical

Week 2 – Causality and equivalent models

- Tuesday May 15 – Lecture + **deadline assignment 1**
- Thursday May 17 - Practical

Week 3 – Time-series analysis and network models

- Tuesday May 22 – Lecture + **deadline assignment 2**
- Thursday May 24 - Free time to work on final project

Week 4 – Wrap-up and presentations

- Tuesday May 29 – Recap + practical + (presentations) + **deadline assignment 3**
- Thursday May 31 - Presentations
- Friday June 1 - **deadline final project report**

Individual Assignments

Each week, the assignment will be made available 15:00 on Tuesday, and will be due 15:00 the next Tuesday. Each assignment will contribute to 20% or 10% (last week) of your grade.

- Work on the assignments **alone**.
- Hand in a **PDF** file and an **.R** file (in case R was used). If you use Jasp, hand in the Jasp object as well as a screenshot of the options used.
- Make sure your PDF report is as standalone readable as possible. E.g., if you are asked to report a factor loading matrix, then report it in the PDF and not just say “look at .R file”.
- Assignments are due **before 15:00**. If you do not hand in an assignment before 15:00, you will get a 1.
- If you encounter any problems, or have any feedback, please let me know before the deadline, as then I can take it into account or help you.

Final Project

Three options:

1. Perform a SEM analysis on your own data and write a report (individual)
2. Write a manual for semPlot, Onyx, Jasp or Lavaan (individual or with a partner)
3. Research an area or a topic of SEM in more detail and teach fellow students about it

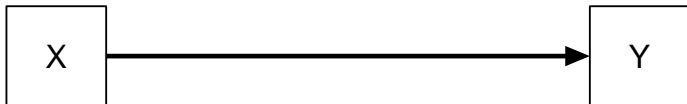
See syllabus on blackboard

- Claim your project using the discussion board on blackboard as soon as possible!
- If you have another idea on a project not listed above, talk to me

Causal modeling

- This course will introduce structural equation modeling (SEM)
- In SEM, we will discuss modeling complex causal hypotheses
- Again, all variables are assumed normally distributed and all associations are assumed linear
- Causal hypotheses can be specified between observed and latent variables
- CFA is a special case of SEM

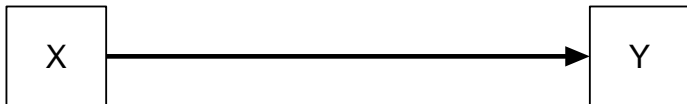
Causal models



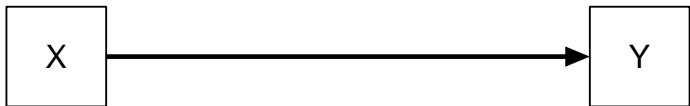
X causes Y

Endogenous and exogenous

- **Exogenous** (independent) variables are variables of which the causal origin are not modeled
 - Exogenous variables have a variance (sometimes not drawn)
 - Exogenous variables, except residuals, are allowed to covary (sometimes not drawn)
 - Latents: ξ (ξ_i); observed: x (x is also used for indicators of latent exogenous variables)
 - Residuals are exogenous
- **Endogenous** (dependent) variables are variables of which the causal origin are modeled
 - Simply stated: endogenous variables have incoming arrows
 - Endogenous variables do **not** have a variance by themselves
 - Latents: η (η); observed: y
 - The causal equation for endogenous variables can be derived from the path diagram by summing all incoming edges

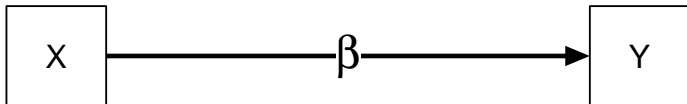


X is exogenous, Y is endogenous

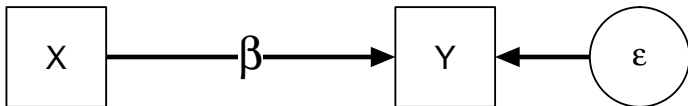


$$y_i = x_i$$

Causal effect goes from right hand side to left hand side.
Experimentally changing x will change y , experimentally changing y will **not** change x

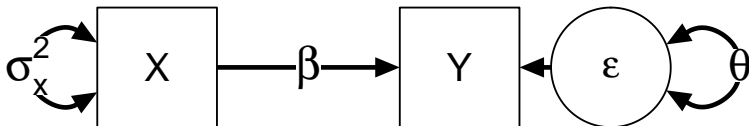


$$y_i = \beta x_i$$



$$y_i = \beta x_i + \varepsilon_i$$

Exogenous variables have a variance (often not drawn)



$$y_i = \beta x_i + \varepsilon_i$$

$$x \sim N(\mu_x, \sigma_x)$$

$$\varepsilon \sim N(0, \theta)$$

$$y_i = \beta x_i + \varepsilon_i$$

$$x \sim N(\mu_x, \sigma_x)$$

$$\varepsilon \sim N(0, \theta)$$

Three observations (variance of x and y and covariance between x and y), three unknowns. Solvable!

$$\text{Var}(x) = \sigma_x^2$$

$$\text{Var}(y) = \beta^2 \sigma_x^2 + \theta$$

$$\text{Cov}(x, y) = \beta \sigma_x^2$$

In addition, we can derive that the *expected value* (mean) of y becomes:

$$\mathcal{E}(y) = \beta \mathcal{E}(x).$$

But why?

Expected Values

Suppose someone offers you to play the following game: You throw a dice, if you throw a 6 you win 6 Euro, but if you throw any other number you pay 1 Euro. Would it, statistically, be smart to play this game?

Expected Values

Suppose someone offers you to play the following game: You throw a dice, if you throw a 6 you win 6 Euro, but if you throw any other number you pay 1 Euro. Would it, statistically, be smart to play this game?

Given a $1/6$ probability to win 6 Euro and a $5/6$ probability to lose 1 Euro, we expect to win:

$$6 \times \frac{1}{6} - 1 \times \frac{5}{6} = 0.166\dots$$

We should play this game!

Expected Values

Let $\mathcal{E}(x)$ denote the *expected value* (mean) of a random variable x that can take K different states (e.g, $K = 2$ when you can win or lose). We can obtain the expected value as follows:

$$\mathcal{E}(x) = \sum_{k=1}^K x_k \Pr(x = x_k).$$



If you bet on a number and the ball falls on that number, you win 35 times your bet! Should you play this game?



If you bet on a number and the ball falls on that number, you win 35 times your bet! Should you play this game?

Suppose we bet 1 Euro, the expected value is:

$$\frac{1}{37} \times 35 - \frac{36}{37} \times 1 = -0.027$$

So you expect to lose 2.7 cent on average for every bet. This is true for *every possible bet* you can place!

Expected Values

If x is continuous, $k \rightarrow \infty$ and $\Pr(x = x_k) \rightarrow 0$, but we can still compute the expected value by using the density function $f(x)$ and integration:

$$\mathcal{E}(x) = \int_{\mathbb{R}} xf(x) dx.$$

More general, for any function $g(x)$, we can obtain:

$$\mathcal{E}(g(x)) = \int_{\mathbb{R}} g(x)f(x) dx.$$

With this we can e.g. proof that if $x \sim N(\mu, \sigma_x^2)$:

$$\begin{aligned}\mathcal{E}(x) &= \mu \\ \text{Var}(x) &= \mathcal{E}((x - \mu)^2) = \sigma_x^2\end{aligned}$$

For example, see: <https://www.statlect.com/probability-distributions/normal-distribution>

Expected Values

From:

$$\mathcal{E}(g(x)) = \int_{\mathbb{R}} g(x)f(x) dx.$$

we can derive the following rules:

$$\mathcal{E}(\alpha) = \alpha$$

$$\mathcal{E}(\alpha x) = \alpha \mathcal{E}(x)$$

$$\mathcal{E}(\alpha x + \beta) = \alpha \mathcal{E}(x) + \beta$$

$$\mathcal{E}(x + y) = \mathcal{E}(x) + \mathcal{E}(y)$$

$$\mathcal{E}(\alpha x + \beta y) = \alpha \mathcal{E}(x) + \beta \mathcal{E}(y)$$

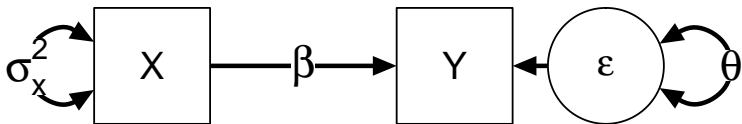
Where α and β are constants (parameter) and x and y are random variables.

Example, the *sample mean* is defined as:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

we can now derive that, if $y \sim N(\mu, \sigma_x^2)$, this is a good estimate for μ :

$$\begin{aligned} \mathcal{E}(\bar{x}) &= \mathcal{E}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \\ &= \frac{1}{n} \mathcal{E}\left(\sum_{i=1}^n x_i\right) \\ &= \frac{1}{n} n \mathcal{E}(x) \\ &= \mu \end{aligned}$$



$$y_i = \beta x_i + \varepsilon_i; x \sim N(\mu_x, \sigma_x); \varepsilon \sim N(0, \theta)$$

We can now derive:

$$\begin{aligned}\mathcal{E}(y) &= \mathcal{E}(\beta x + \varepsilon) \\ &= \mathcal{E}(\beta x) + \mathcal{E}(\varepsilon) \\ &= \beta \mathcal{E}(x)\end{aligned}$$

Expected Values

Multivariate generalizations are straightforward:

$$\mathcal{E}(\mathbf{A}) = \mathbf{A}$$

$$\mathcal{E}(\mathbf{Ax}) = \mathbf{A}\mathcal{E}(\mathbf{x})$$

$$\mathcal{E}(\mathbf{xB}) = \mathcal{E}(\mathbf{x})\mathbf{B}$$

$$\mathcal{E}(\mathbf{Ax} + \mathbf{B}) = \mathbf{A}\mathcal{E}(\mathbf{x}) + \mathbf{B}$$

$$\mathcal{E}(\mathbf{x} + \mathbf{y}) = \mathcal{E}(\mathbf{x}) + \mathcal{E}(\mathbf{y})$$

Where \mathbf{A} and \mathbf{B} are constant (parameter) matrices. For example, given $\mathcal{E}(\boldsymbol{\eta}) = \boldsymbol{\alpha}$:

$$\mathbf{y}_i = \lambda\boldsymbol{\eta}_i + \boldsymbol{\varepsilon}_i$$

$$\mathcal{E}(\mathbf{y}) = \mathcal{E}(\lambda\boldsymbol{\eta} + \boldsymbol{\varepsilon})$$

$$\mathcal{E}(\mathbf{y}) = \lambda\mathcal{E}(\boldsymbol{\eta}) + \mathcal{E}(\boldsymbol{\varepsilon})$$

$$\mathcal{E}(\mathbf{y}) = \lambda\boldsymbol{\alpha}$$



Given that $\text{Var}(x) = \mathcal{E}((x - \mu_x)^2)$ and $\mathcal{E}(x) = \mu_x = -0.027$ for betting 1 Euro in Roulette. What is the variance and standard deviation of our bet?



Given that $\text{Var}(x) = \mathcal{E}((x - \mu_x)^2)$ and $\mathcal{E}(x) = \mu_x = -0.027$ for betting 1 Euro in Roulette. What is the variance and standard deviation of our bet?

$$\begin{aligned} \text{Var}(\text{bet 1 Euro}) &= \frac{1}{37} \times (35 - -0.027)^2 + \frac{36}{37} \times (1 - -0.027)^2 \\ &= 34.19 \end{aligned}$$

$$\text{SD}(\text{bet 1 Euro}) = \sqrt{\text{Var}(\text{bet 1 Euro})} = 5.85$$

Covariance Algebra

Let $\text{Var}(x)$ indicate “the variance of x ” and $\text{Cov}(x, y)$ indicate “the covariance between x and y ”. Given that $\text{Cov}(x, y) = \mathcal{E}((x - \mu_x)(y - \mu_y))$, the following rules can be derived:

$$\text{Var}(x) = \text{Cov}(x, x)$$

$$\text{Cov}(x, \alpha) = 0$$

$$\text{Cov}(x, y) = \text{Cov}(y, x)$$

$$\text{Cov}(\alpha x, \beta y) = \alpha\beta \text{Cov}(x, y)$$

$$\text{Cov}(x + y, z) = \text{Cov}(x, z) + \text{Cov}(y, z)$$

Where α and β are constants (parameter) and x , y , and z are random variables.

Covariance Algebra

Some consequences:

$$\begin{aligned}\text{Cov}(\alpha x + \beta y, z) &= \text{Cov}(\alpha x, z) + \text{Cov}(\beta y, z) \\ &= \alpha \text{Cov}(x, z) + \beta \text{Cov}(y, z)\end{aligned}$$

$$\text{Var}(x + y) = \text{Var}(x) + \text{Var}(y) + 2\text{Cov}(x, y)$$

$$\text{Var}(\beta x) = \beta^2 \text{Var}(x)$$

Where α and β are constants (parameter) and x , y , and z are random variables.

Matrix Covariance Algebra

Let $\text{Var}(\mathbf{x})$ indicate “the variance–covariance matrix of vector \mathbf{x} ” and $\text{Cov}(\mathbf{x}, \mathbf{y})$ indicate “the covariance matrix between \mathbf{x} and \mathbf{y} ”. Then the following rules can be derived:

$$\text{Var}(\mathbf{x}) = \text{Cov}(\mathbf{x}, \mathbf{x})$$

$$\text{Cov}(\mathbf{Ax}, \mathbf{By}) = \mathbf{ACov}(\mathbf{x}, \mathbf{y})\mathbf{B}^\top$$

$$\text{Var}(\mathbf{Bx}) = \mathbf{BVar}(\mathbf{x})\mathbf{B}^\top$$

$$\text{Cov}(\mathbf{x} + \mathbf{y}, \mathbf{z}) = \text{Cov}(\mathbf{x}, \mathbf{z}) + \text{Cov}(\mathbf{y}, \mathbf{z})$$

Where \mathbf{A} and \mathbf{B} are constant (parameter) matrices.

$$y_i = \beta x_i + \varepsilon_i$$

$$\text{Var}(x) = \sigma_x^2$$

$$\begin{aligned}\text{Var}(y) &= \text{Var}(\beta x + \varepsilon) \\ &= \text{Cov}(\beta x + \varepsilon, \beta x + \varepsilon) \\ &= \text{Cov}(\beta x, \beta x + \varepsilon) + \text{Cov}(\varepsilon, \beta x + \varepsilon) \\ &= \text{Cov}(\beta x, \beta x) + \text{Cov}(\beta x, \varepsilon) + \text{Cov}(\varepsilon, \beta x) + \text{Cov}(\varepsilon, \varepsilon)\end{aligned}$$

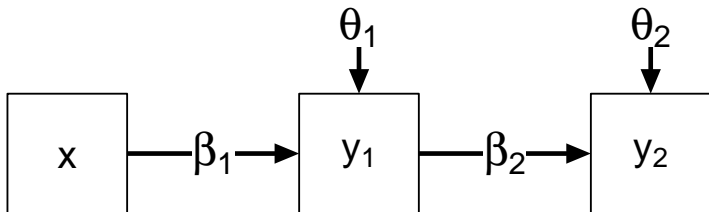
But since x is not correlated with the residuals, $\text{Cov}(x, \varepsilon) = 0$ and thus:

$$\begin{aligned}\text{Var}(y) &= \beta^2 \text{Cov}(x, x) + \text{Cov}(\varepsilon, \varepsilon) \\ &= \beta^2 \text{Var}(x) + \text{Var}(\varepsilon)\end{aligned}$$

$$y_i = \beta x_i + \varepsilon_i$$

$$\begin{aligned}\text{Cov}(x, y) &= \text{Cov}(x, \beta x_i + \varepsilon_i) \\ &= \text{Cov}(x, \beta x_i) + \text{Cov}(x, \varepsilon_i) \\ &= \beta \text{Cov}(x, x_i) \\ &= \beta \text{Var}(x)\end{aligned}$$

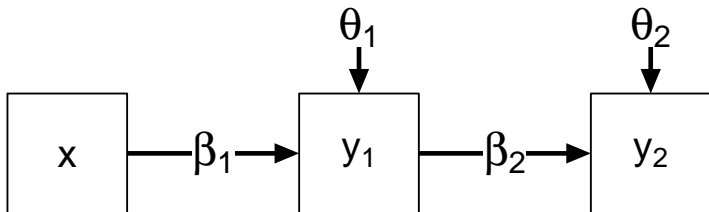
Path analysis



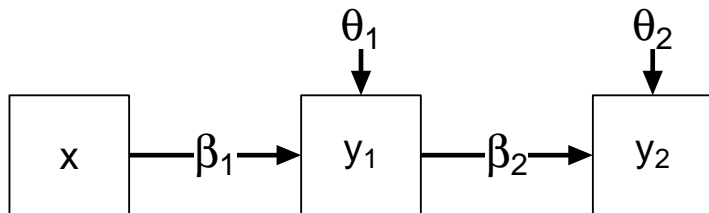
x is exogenous, and both y_1 and y_2 are endogenous. θ_1 is the variance of ε_1 . Causal model for y_2 :

$$y_{i2} = \beta_2 y_{i1} + \varepsilon_{i2}$$

$$y_{i2} = \beta_2(\beta_1 x_i + \varepsilon_{i1}) + \varepsilon_{i2}$$



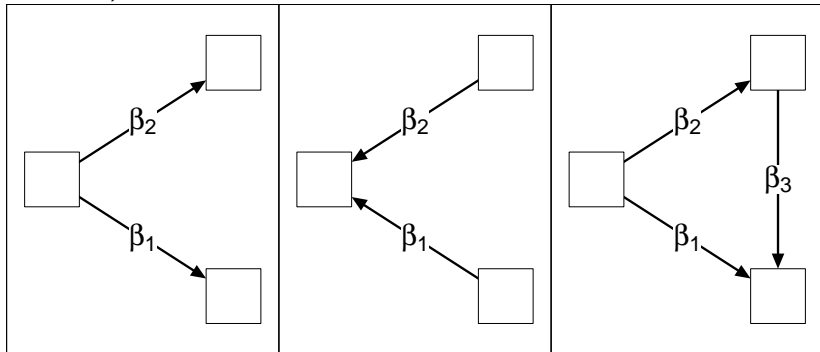
Number of parameters: 2 regressions +2 residual variances +1 exogenous variance (not drawn) = 5, number of observations: 3 variances and 3 covariances. 1 degree of freedom!



Implied covariance between x and y_2 :

$$\begin{aligned}\text{Cov}(x, y_2) &= \text{Cov}(x, \beta_2(\beta_1 x_i + \varepsilon_{i1}) + \varepsilon_{i2}) \\ &= \text{Cov}(x, \beta_2 \beta_1 x + \beta_2 \varepsilon_1 + \varepsilon_2) \\ &= \text{Cov}(x, \beta_2 \beta_1 x) + \text{Cov}(x, \beta_2 \varepsilon_1) + \text{Cov}(x, \varepsilon_2) \\ &= \beta_1 \beta_2 \text{Cov}(x, x) \\ &= \beta_1 \beta_2 \sigma_x\end{aligned}$$

Practical: identify exogenous (x) and endogenous (y) variables, and derive expressions for expected values and (co)variances of/between all variables, in terms of parameters and expectations/(co)variances of exogenous variables (x and residuals):



Also think about the final project!