

# Modèles d'équations structureaux

## Contents

|   |          |
|---|----------|
| <b>Introduction</b>   | <b>1</b> |
| <b>Contenu du cours</b>   | <b>1</b> |
| <b>Types of variables and relationships in a structural equation model:</b> | <b>1</b> |
| From theoretical model to statistical model: . . . . .                      | 5        |
| <b>Model Fitting in Lavaan</b>  | <b>8</b> |
| References . . . . .  | 10       |

## Introduction

Structural equation models belong to a family of models consisting of a set of mathematical equations and assumptions about a studied system. These assumptions stem from our prior knowledge or assumptions about how the system operates. In statistics, a system is a set of variables or phenomena that are studied within the framework of an analysis or study. These variables may be related by complex relationships. Our aim with the analysis is to understand how these variables interact with each other or how they influence a specific outcome or phenomenon.

## Contenu du cours

- Types de variables et relations dans un modèle d'équations structurelles ;
- Du modèle théorique au modèle statistique ;
- Ajustement du modèle dans lavaan.

## Types of variables and relationships in a structural equation model:

The theoretical structure of a structural equation model encompasses several types of variables, defining their characteristics and roles within the model.

Based on their nature, variables can be classified into 1) latent variables and 2) observed variables. A latent variable is a variable that is not directly measured but represents concepts or traits that do not have a clear unit of measurement. This type of variable is often used in psychology (intelligence, satisfaction, etc.). Observed variables are measured or collected using established methods within the discipline. In ecology, we often work with observed variables.

Based on the role of variables in the theoretical model, variables can be exogenous or endogenous. Exogenous variables are independent variables that influence other variables in the model but are not influenced by any other variable in return. They represent the drivers of changes in our system. Endogenous variables, on the other hand, are variables influenced by exogenous variables or other variables in the model, and typically represent the core of our system and the outcomes of the processes we are describing with our model.

Depending on how variables are conceptualized in the model, variables may have moderator and mediator status. A moderator is a variable that influences the strength and direction of the relationship between two variables. A moderator does not explain the “causes” of the relationship but only intervenes in the quantitative aspects of the relationship between two variables. A mediator is a variable that explains the relationship between an independent variable and a dependent variable. A variable acts as a mediator when it represents the causal mechanism linking two variables.

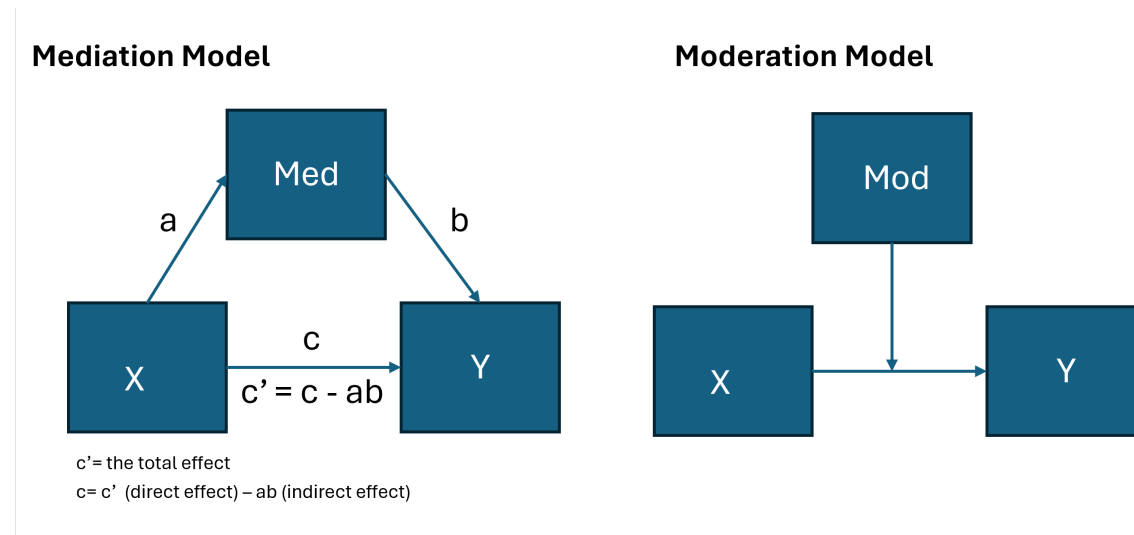


Figure 1: Models for mediation analysis and moderator analysis

In the figure, the arrows represent paths. In a structural equation model, a path represents a relationship that connects two variables, which can be either direct or indirect, and may include mediated effects. As you can see in the image, mediated effects refer to a situation where the impact of an independent variable on a dependent variable passes through another variable. In other words, the mediating variable transmits or mediates the effect of the independent variable on the dependent variable.

In a structural equation model, mediation and moderation analysis are integrated to gain a deeper understanding of the relationships between our variables, considering both mediating processes and the effects of moderators. However, mediation and moderation analyses can be conducted separately if we want to test relationships among three variables using the Psych package:

The ‘mediate’ function in the ‘psych’ package allows you to conduct a mediation analysis with the ‘mediate’ function. The mediating variable must be enclosed in parentheses to inform the function of its role in the model. In the example we will use to illustrate the function, we use the ‘mediate’ function to test the direct and indirect effects of temperature on the width of growth rings. Indeed, in the northern hemisphere and in cold environments where temperature is a limiting factor for growth, the width of growth rings increases with temperature. However, the width of the growth ring is closely related to the number of wood cells that make up the ring, which is, in turn, related to temperature. Thus, a portion of the total effect of temperature on the width of the ring is indeed mediated by the number of cells. If the direct effect of temperature on the number of cells is greater than that on the width of the ring, temperature changes will be more related to changes in the number of cells than to the width of the ring, and therefore the latter will be more easily predictable.

```
require(psych)
```

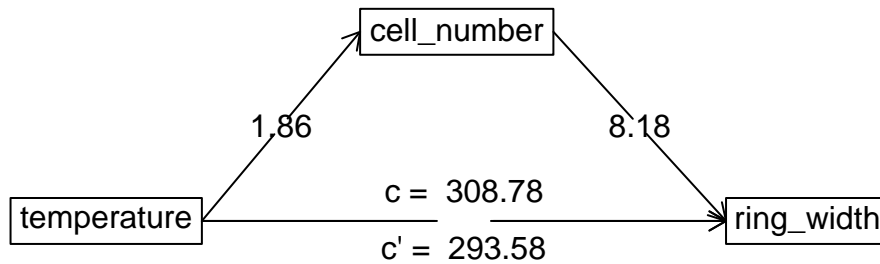
```
## Le chargement a nécessité le package : psych
```

```
## Warning: le package 'psych' a été compilé avec la version R 4.3.3
```

```
ringdatacell <- read.csv("C:/Users/buttoval/Documents/ECL8202/donnees/ringdatacell.CSV")

medanalysis<-mediate( ring_width ~ temperature + (cell_number) , data=ringdatacell )
```

## Mediation



“In the diagram, it is clear that the total effect of temperature is  $c = 8.15$ , but the direct effect of temperature on growth ring is much smaller,  $c' = 0.56$ .

Nota bene: The ‘mediate’ function provides standardized coefficients and centered means by default. To avoid this, use the ‘zero=FALSE’ and ‘std=FALSE’ arguments.”

```
summary(medanalysis)
```

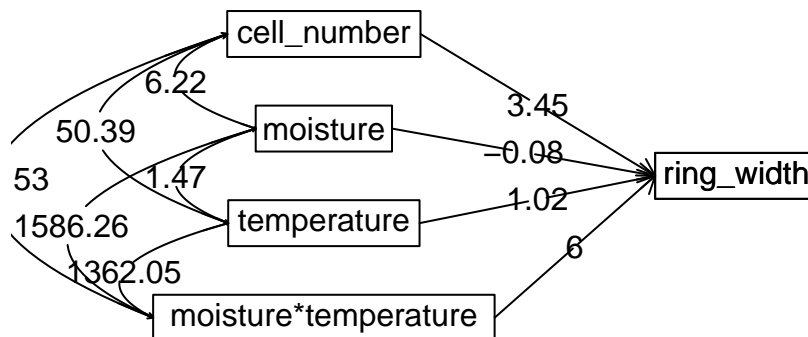
```
## Call: mediate(y = ring_width ~ temperature + (cell_number), data = ringdatacell)
##
## Direct effect estimates (traditional regression)      (c') X + M on Y
##      ring_width      se      t df      Prob
## Intercept      -157.63 459.43 -0.34 97 7.32e-01
## temperature      293.58 44.61  6.58 97 2.38e-09
## cell_number       8.18 21.29  0.38 97 7.02e-01
##
## R = 0.84 R2 = 0.7    F = 112.67 on 2 and 97 DF    p-value: 5.08e-26
##
## Total effect estimates (c) (X on Y)
##      ring_width      se      t df      Prob
## Intercept      -93.82 426.50 -0.22 98 8.26e-01
## temperature      308.78 20.49 15.07 98 2.88e-27
```

```
##
## 'a' effect estimates (X on M)
##           cell_number  se    t df    Prob
## Intercept           7.80 2.03  3.84 98 2.20e-04
## temperature         1.86 0.10 19.04 98 1.06e-34
##
## 'b' effect estimates (M on Y controlling for X)
##           ring_width    se    t df    Prob
## cell_number           8.18 21.29 0.38 97 0.702
##
## 'ab' effect estimates (through all mediators)
##           ring_width boot    sd lower upper
## temperature        15.21 13.84 42.95 -70.96 96.26
```

In a moderation analysis, we examine how the effect of an independent variable on a dependent variable may be modified by another variable, known as a moderator. To assess this interaction, we typically include an interaction term in the multiple regression model. This interaction term allows us to test whether the effect of the independent variable on the dependent variable varies depending on the levels of the moderator. In summary, a moderation analysis is a multiple regression with an interaction. We can obtain a moderation model using the 'mediate' function, but without specifying the effect of the moderator

```
medanalysisys<-mediate( ring_width ~ cell_number + moisture + moisture*temperature,
                        data=ringdatacell,zero= FALSE,std=FALSE)
```

## Moderation model



```
summary(medanalysisys)
```

```
## Call: mediate(y = ring_width ~ cell_number + moisture + moisture *
##             temperature, data = ringdatacell, std = FALSE, zero = FALSE)
```

```
##
## No mediator specified leads to traditional regression
##           ring_width    se      t df      Prob
## Intercept           9.04 11.16   0.81 95  4.20e-01
## cell_number         3.45  0.13  27.50 95  4.97e-47
## moisture          -0.08  0.22  -0.33 95  7.39e-01
## temperature         1.02  0.66   1.55 95  1.25e-01
## moisture*temperature  6.00  0.01 522.12 95  4.55e-166
##
## R = 1 R2 = 1    F = 2400995 on 4 and 95 DF    p-value:  9.14e-237
```

## From theoretical model to statistical model:

When deciding to conduct a SEM, a well-defined hypothesis represents the best investment to leverage this analysis. Therefore, we will classify our variables according to the typology defined in the previous paragraph, and then construct our a priori model. This model represents our understanding of the system based on the scientific evidence we have gathered in our study and contains our hypotheses in the form of links between variables.

In a SEM, we establish a theoretical model based on postulated relationships between variables, and then test this model with real data to see if it fits these data well. The objective is to determine if the theoretical model is statistically valid and can be generalized to real data. Model validation thus requires not rejecting the null hypothesis that the relationships between variables as specified in the theoretical model are also present in the real data, and that any observed difference between the theoretical model and real data is due to chance or measurement errors.

In a SEM, symbols can be used to represent links and variables in a diagram:



Figure 2: Symbols and diagrams

There are several packages available in R that allow you to perform a SEM. Here, we will use Lavaan, which has its own syntax for defining the variables of the model and their links.

| Formula and Definition | Operator   | Meaning            |
|------------------------|------------|--------------------|
| Latent variable        | $\sim$     | is obtained from   |
| Covariate              | $\sim\sim$ | is correlated with |
| intercept              | $\sim 1$   | intercept          |

For an example, we will use a simulated dataset containing the following information:

Temperature: The average temperature in degrees Celsius recorded during the growing season. Humidity: The average percentage of relative humidity recorded during the growing season. Stem Size: The average stem size of the plant, in cm. Cell Number: The total number of cells observed in each growth ring. Ring Width: The average width of growth rings, a measure of annual tree growth.

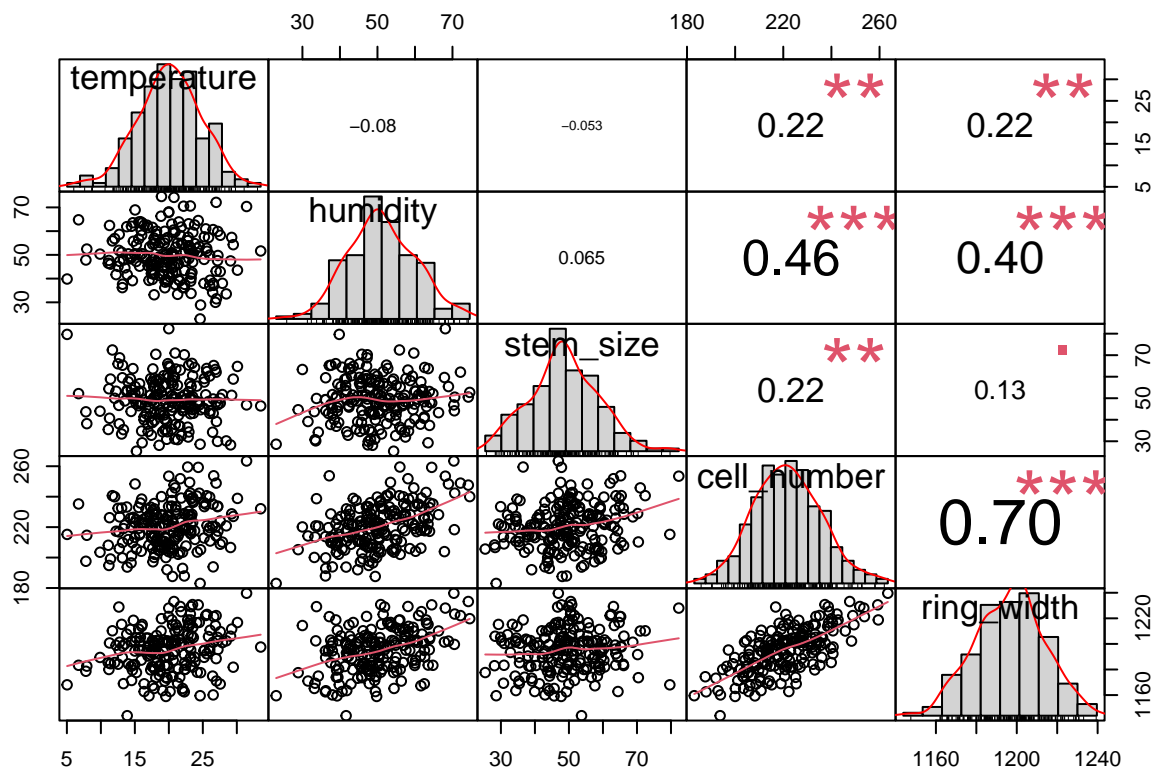
```
simulated_data <- read.csv("C:/Users/buttoval/Documents/ECL8202/donnees/simulatedsemring.csv")
```

Before adjusting a SEM model it is often useful to visualize a correlation matrix between the variables:

```
library(PerformanceAnalytics)
```

```
## Le chargement a nécessité le package : xts
```

```
## Le chargement a nécessité le package : zoo
##
## Attachement du package : 'zoo'
## Les objets suivants sont masqués depuis 'package:base':
##
##   as.Date, as.Date.numeric
##
## Attachement du package : 'PerformanceAnalytics'
## L'objet suivant est masqué depuis 'package:graphics':
##
##   legend
chart.Correlation(simulated_data, histogram = TRUE, method = "pearson")
```



The correlation matrix displays high correlations between most of our variables, but it does not tell us anything about the relationships among them.

A theoretical model based on the literature is proposed to explain the relationships between our variables and the underlying process of growth:

On lavaan, we can translate the model using the following syntax:

```
require(lavaan)
```

```
## Le chargement a nécessité le package : lavaan
```

```
## This is lavaan 0.6-15
```

## Sem structure – Theoretical model

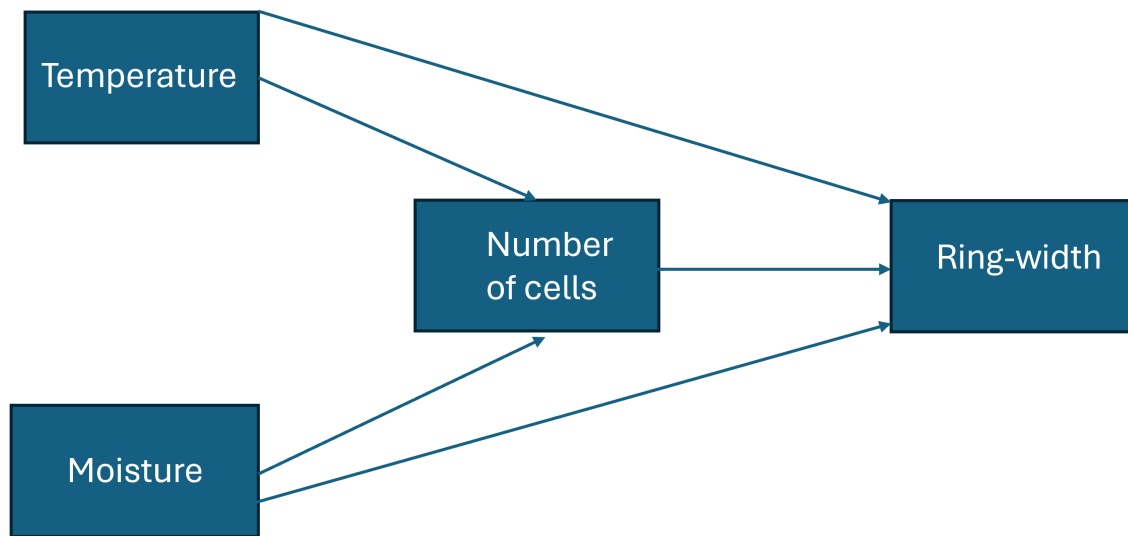


Figure 3: Relationship between environmental factors and growth: each link is supported by literature sources

```
## lavaan is FREE software! Please report any bugs.
##
## Attachement du package : 'lavaan'
## L'objet suivant est masqué depuis 'package:psych':
##
##      cor2cov
```

```
myModel <- '
# regressions
ring_width ~ temperature + humidity+ cell_number
cell_number ~ temperature + humidity
cell_number ~ ~ stem_size
'
```

```
fit <- sem(model = myModel,
           data = simulated_data)
```

```
summary(fit)
```

```
## lavaan 0.6.15 ended normally after 28 iterations
##
##      Estimator              ML
##      Optimization method    NLMINB
##      Number of model parameters      9
##
##      Number of observations      200
```

```
##
## Model Test User Model:
##
##   Test statistic                1.368
##   Degrees of freedom            3
##   P-value (Chi-square)          0.713
##
## Parameter Estimates:
##
##   Standard errors                Standard
##   Information                    Expected
##   Information saturated (h1) model Structured
##
## Regressions:
##           Estimate Std.Err z-value P(>|z|)
##   ring_width ~
##     temperature      0.343   0.187   1.832   0.067
##     humidity         0.235   0.105   2.239   0.025
##     cell_number      0.735   0.069  10.592   0.000
##   cell_number ~
##     temperature      0.803   0.177   4.529   0.000
##     humidity         0.724   0.091   7.934   0.000
##
## Covariances:
##           Estimate Std.Err z-value P(>|z|)
##   .cell_number ~~
##     stem_size        30.682   9.433   3.253   0.001
##
## Variances:
##           Estimate Std.Err z-value P(>|z|)
##   .ring_width       150.351  15.035  10.000   0.000
##   .cell_number       156.311  15.631  10.000   0.000
##   stem_size         107.827  10.783  10.000   0.000
```

## Model Fitting in Lavaan

We consider a SEM to adequately represent our data distribution and natural phenomenon when the p-value (chi-square) of the model is not significant. This is because we aim to observe a correspondence between the observed data and the predictions of our model, which should not exhibit significant differences. In this case, the p-value is 0.713, which reassures us about the model's ability to represent our phenomenon. It is also useful to check other model fit measures for a more comprehensive evaluation of model adequacy. Here are the commonly used indicators, as summarized by Joreskog, K., & Sorbom, D. (1993):

Il existe des standards pour nous naviguer dans la présentation de ces indicateurs:

The summary of the sem function provides us with regression coefficients that indicate the strength and direction of the relationship between our variables. Covariance coefficients also measure the strength and direction of the correlation between variables. These coefficients can be interpreted like any coefficients in a linear regression. However, it can be very useful to standardize the coefficients if we want to compare the effect of different variables on our response variable.

```
summary(fit,standardized=TRUE)
```

```
## lavaan 0.6.15 ended normally after 28 iterations
##
```



**Table 1: Fit indices and their acceptable thresholds**

| Fit Index  | Acceptable Threshold Levels  | Description   |
|--|--|---|
| <i>Absolute Fit Indices</i><br>Chi-Square $\chi^2$ | Low $\chi^2$ relative to degrees of freedom with an insignificant $p$ value ( $p > 0.05$ ) |   |
| Relative $\chi^2$ ( $\chi^2/\text{df}$ )           | 2:1 (Tabachnik and Fidell, 2007)<br>3:1 (Kline, 2005)                                      | Adjusts for sample size.  |
| Root Mean Square Error of Approximation (RMSEA)    | Values less than 0.07 (Steiger, 2007)  | Has a known distribution. Favours parsimony. Values less than 0.03 represent excellent fit.   |
| GFI  | Values greater than 0.95   | Scaled between 0 and 1, with higher values indicating better model fit. This statistic should be used with caution.   |
| AGFI   | Values greater than 0.95   | Adjusts the GFI based on the number of parameters in the model. Values can fall outside the 0-1.0 range.  |
| RMR  | Good models have small RMR (Tabachnik and Fidell, 2007)                                    | Residual based. The average squared differences between the residuals of the sample covariances and the residuals of the estimated covariances. Unstandardised. |
| SRMR   | SRMR less than 0.08 (Hu and Bentler, 1999)   | Standardised version of the RMR. Easier to interpret due to its standardised nature.  |
| <i>Incremental Fit Indices</i><br>NFI              | Values greater than 0.95   | Assesses fit relative to a baseline model which assumes no covariances between the observed variables. Has a tendency to overestimate fit in small samples.     |
| NNFI (TLI)   | Values greater than 0.95   | Non-normed, values can fall outside the 0-1 range. Favours parsimony. Performs well in simulation studies (Sharma et al, 2005; McDonald and Marsh, 1990)        |
| CFI  | Values greater than 0.95   | Normed, 0-1 range.  |

Figure 4: Model adjustment

**Table 2: Hu and Bentler's Two-Index Presentation Strategy (1999)**

| Fit Index Combination | Combinational Rules                                |
|-----------------------|--|
| NNFI (TLI) and SRMR   | NNFI of 0.96 or higher and an SRMR of .09 or lower |
| RMSEA and SRMR        | RMSEA of 0.06 or lower and a SRMR of 0.09 or lower |
| CFI and SRMR          | CFI of .96 or higher and a SRMR of 0.09 or lower   |

Figure 5: Acceptable combinations of diagnostic indicators

```

##      Estimator                      ML
##      Optimization method          NLMINB
##      Number of model parameters      9
##
##      Number of observations          200
##
## Model Test User Model:
##
##      Test statistic                  1.368
##      Degrees of freedom              3
##      P-value (Chi-square)            0.713
##
## Parameter Estimates:
##
##      Standard errors                Standard
##      Information                    Expected
##      Information saturated (h1) model Structured
##
## Regressions:
##      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##      ring_width ~
##      temperature      0.343   0.187   1.832   0.067   0.343   0.096
##      humidity          0.235   0.105   2.239   0.025   0.235   0.128
##      cell_number       0.735   0.069  10.592   0.000   0.735   0.618
##      cell_number ~
##      temperature      0.803   0.177   4.529   0.000   0.803   0.267
##      humidity          0.724   0.091   7.934   0.000   0.724   0.467
##
## Covariances:
##      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##      .cell_number ~~
##      stem_size        30.682   9.433   3.253   0.001  30.682   0.236
##
## Variances:
##      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
##      .ring_width      150.351  15.035  10.000   0.000  150.351   0.497
##      .cell_number     156.311  15.631  10.000   0.000  156.311   0.730
##      stem_size        107.827  10.783  10.000   0.000  107.827   1.000

```

## References

- Revelle, 2024: How to use the psych package for regression and mediation analysis: <https://cran.r-project.org/web/packages/psychTools/vignettes/mediation.pdf>
- Yves Rosseel (2012). lavaan: An R Package for Structural Equation Modeling. Journal of Statistical Software, 48(2), 1-36. URL <http://www.jstatsoft.org/v48/i02/> <https://lavaan.ugent.be/>
- Joreskog, K., & Sorbom, D. (1993). Structural equation modelling: Guidelines for determining model fit. NY: University Press of America.