

Twin Structural Equation Models in R & Lavaan

Michel G Nivard

6/21/2021

Contents

1	Introduction	2
1.1	Resemblance between relatives	2
1.2	SEM Models for twin data	2
2	Quickly set up my R environment	3
3	Lavaan	3
3.1	Syntax	3
3.2	Estimators	5
4	Single trait twin model in lavaan	5
4.1	ACE model	5
4.2	ADE model	10
4.3	Sibling interactions	14
4.4	Binary/Ordinal data	20
4.4.1	Ordinal data	20
4.4.2	Binary data	25
4.5	Sex specific effects	30
4.6	Gene-environment interaction	30
4.7	Gene-environment correlation	31
4.8	Rater bias models	31
5	Multitrait twin models in lavaan	31
5.1	Cholesky decomposition	31
5.2	Direction of Causation models	31
6	Longitudinal twin models	31
6.1	Auto regressive models	31
6.1.1	Phenotype to environment effects	31
6.2	Growth curve models	31

1 Introduction

1.1 Resemblance between relatives

There is a long history of questioning the nature of the resemblance between family members, specifically sibling and even more specifically twins. Reasons that are commonly put forward for these similarities are the obvious genetic similarity, similarity in socio-economic position and similarity in upbringing shared between siblings. Twin and family models leverage variation in genetic and environmental relatedness between family members to estimate the relative contributions of genetics, the environment, their interacting, their correlation and other process to the similarities between relatives.

Twins offer an excellent “natural experiment” where some twins are genetically identical (identical, or monozygotic twins or MZ twins or MZs) and some twins share half their segregating DNA (fraternal, or dizygotic twins, DZ twins or DZs). Like any natural experiment, twin studies aren’t true experiments, and so the estimated quantities come with various assumptions, some of which are specific to twin models, and I’ll make sure to catalog the assumptions we are making along the way.

If certain assumptions are satisfied the causes of resemblance between twins can be generalized to be causes on individual differences in the general population, provided the twins we sample are representative of the population, and the population in question is homogeneous, as the relative contribution of genes, the environment their correlation and interaction aren’t fixed quantities or inescapable facts of life. these estimates are temporal phenomena that reflect the economic/political and societal status quo in the sample of population one studies.

1.2 SEM Models for twin data

This document describes various twin structural equation models (SEM) to estimate the nature of the relationship between family members in order to learn about the contribution of genes and the social or rearing environment to complex (behavioral) outcomes. The goal is to provide a basic understands of these models with the means to fit the models in lavaan (Rosseel 2012), lavaan is an R package that allow the use to define a structural equation model in terms of regression, variances and covariances and will be familiar to users of M-Plus. The package is (IMO) more accessible to beginners then another excellent SEM R package: OpenMx (Neale et al. 2015), but the accessibility comes at the cost of less flexibility, In terms of flexibility OpenMx is truly unrivaled. Its worth pointing out that the developers behind OpenMx have an academic interest in twin models, which means scripts and support for users in those models is often especially excellent. The best (to my knowledge) repository of twin models in OpenMx is maintained by Hermine Maes and there is an R package that wraps OpenMx twin models in simpler R commends called umx (Bates, Maes, and Neale 2019) which is build and maintained by Tim Bates. Some models considered fundamental to our understanding of gene-environment interaction cannot (to my knowledge) be fitted in lavaan, despite these limitation there are many advantages to learning lavaan. The syntax used to specify lavaan models also features in blavaan (Merkle and Rosseel 2018), its Bayesian cousin providing an easy entry into Bayesian twin modeling. Similarly the package RegSEM (Jacobucci, Grimm, and McArdle 2016) and lslx (Huang, Chen, and Weng 2017) allows for *regularized* structural equation modeling and use lavaan, or very lavaan like, syntax and could be a point of departure for various novel and innovative twin models.

All things considered it makes sense to collect and document lavaan syntax and examples for various twin models.

2 Quickly set up my R environment

We will need the following packages, some may be optional some will be required for any lavaan analysis (like lavaan).

```
library(lavaan)
library(MASS)
library(tidySEM)
library(ggplot2)
library(knitr)
```

3 Lavaan

3.1 Syntax

Lavaan requires the user to specify a SEM model in terms of regression, which are directional relations between observed variables, factor loadings which are directional relations between a latent factor and an observed variable, variances and variances which are undirected relationships between variables. The model is defined in a text string within R.

Table 1: Table 1: Lavaan Syntax Elements and their meaning

Syntax	Description	Meaning
f =~ In1 + In2 + In3	Factor loading	The factor “f” is measure by 3 indicators “In1,” “In2” and “In3,”
y ~ x	Regression	The variable “y” is regressed on the variable “x”
x1 ~~ x2	Covariance	the variable “x1” and “x2” covary
x1 ~~ x1	Variance	the variable x1 has a freely estimated variance

You can use addition to add multiple elements to a regression in lavaan: $y \sim x1 + x2$ and to covariance: $x1 \sim x1 + x2 + x3$. Lavaan will map the variable names in the syntax to variable names in your dataset. Lavaan expects all the relationships in a model to be contained in a single string variable in R let me show you an example model fit to simulated data:

```
# make the data
f <- rnorm(100)
i1 <- 1*f + rnorm(100)
i2 <- 0.6*f + rnorm(100)
i3 <- 1.4*f + rnorm(100)

x <- rnorm(100)
y <- f + x

dataset <- cbind.data.frame(i1,i2,i3,x,y)

# write the lavaan model:
example.model <- "
f1 =~ i1 + i2 + i3
y ~ x + f1
x ~~ f1
```

```
"
```

Okay, now notice how `f` does not exist in the dataframe: `dataset` (its going to be latent), lets fit the model to the data and have a look at the results.

```
model.fit <- sem(model = example.model,data=dataset)
```

```
## Warning in lav_object_post_check(object): lavaan WARNING: some estimated ov  
## variances are negative
```

```
model.fit
```

```
## lavaan 0.6-7 ended normally after 30 iterations  
##  
##      Estimator                      ML  
##      Optimization method          NLMINB  
##      Number of free parameters      11  
##  
##      Number of observations          100  
##  
## Model Test User Model:  
##  
##      Test statistic                  4.888  
##      Degrees of freedom              4  
##      P-value (Chi-square)           0.299
```

So the model ran and converged, the fit (according to the `chi2`) is adequate. Thats not entirely unexpected as we are fitting the true model to the data. Lets look at the estimates:

```
summary(model.fit)
```

```
## lavaan 0.6-7 ended normally after 30 iterations  
##  
##      Estimator                      ML  
##      Optimization method          NLMINB  
##      Number of free parameters      11  
##  
##      Number of observations          100  
##  
## Model Test User Model:  
##  
##      Test statistic                  4.888  
##      Degrees of freedom              4  
##      P-value (Chi-square)           0.299  
##  
## Parameter Estimates:  
##  
##      Standard errors                Standard  
##      Information                    Expected  
##      Information saturated (h1) model Structured  
##
```

```

## Latent Variables:
##           Estimate Std.Err z-value P(>|z|)
##   f1 =~
##     i1           1.000
##     i2           0.472    0.111    4.241    0.000
##     i3           1.397    0.165    8.461    0.000
##
## Regressions:
##           Estimate Std.Err z-value P(>|z|)
##   y ~
##     x           1.040    0.051   20.250    0.000
##     f1           1.043    0.120    8.718    0.000
##
## Covariances:
##           Estimate Std.Err z-value P(>|z|)
##   f1 ~~
##     x          -0.144    0.112   -1.282    0.200
##
## Variances:
##           Estimate Std.Err z-value P(>|z|)
##   .i1           0.951    0.143    6.653    0.000
##   .i2           0.989    0.139    7.088    0.000
##   .i3           0.780    0.147    5.292    0.000
##   .y          -0.032    0.055   -0.579    0.563
##   x            1.021    0.144    7.071    0.000
##   f1            0.936    0.236    3.975    0.000

```

As you can see the factor loading are close to their true simulated values of 1, 0.6 and 1.4 as are all of the other parameters.

3.2 Estimators

Int he examples below is simulate data and have no missingness or non-normal data to deal with, in reality you will so it makes sense to discuss the various estimators you could use in lavaan.

4 Single trait twin model in lavaan

4.1 ACE model

The “ACE” model operationalized a phenotype as a function of additive genetic variance (A), common environmental influences(C) (can be rearing, can be societal influences can be governmental policies), and environment unique to each of the individual twins (E) (can be private friends, being in separate classrooms, but also measurement error). The correlation between the genetic influences and common environmental is 1 for MZ twins while the correlation between the genetic influences is 0.5 for DZ twins.

the ACE model assumes among other things:

1. The absence of non-additive genetic effects (gene gene interaction and dominance)
2. The absence of sibling influences or rater contrast effect
3. The absence of gene by environment correlation
4. The absence of gene by environment interaction

5. The environments MZ twin grow up are as similar as the environments DZ twins grow up in.

Lets simulate same data where the resemblance between twins is a function of equal parts (33.3%/33.3%/33.4%)

```
A <- matrix(1,2,2) # genetic correlation for MZ's = 1
C <- matrix(1,2,2)
E <- diag(2)
Adz <- matrix(c(1,.5,.5,1),2,2) # genetic correlation for DZ's = 0.5

# make 1000 pairs of MZ twins
MZ <- mvrnorm(1000,mu=c(0,0),Sigma = A+C+E)

# Add a column to label as MZ:
MZ<- cbind.data.frame("MZ",MZ)
colnames(MZ) <- c("zyg","P1", "P2")

# make 1500 DZ twin pairs
DZ <- mvrnorm(1500,mu=c(0,0),Sigma = Adz+C+E)

# add variable too label as DZ:
DZ <- cbind.data.frame("DZ",DZ)
colnames(DZ) <- c("zyg","P1", "P2")

# Combine MZ and DZ twins
dataset <- rbind(MZ,DZ)
```

We then define the lavaan model that can express the variance in the trait P explained by latent variables A, C and E:

```
ace.model<-"
A1=~ NA*P1 + c(a,a)*P1
A2=~ NA*P2 + c(a,a)*P2
C1 =~ NA*P1 + c(c,c)*P1
C2 =~ NA*P2 + c(c,c)*P2
# variances
A1 ~~ 1*A1
A2 ~~ 1*A2
C1 ~~ 1*C1
C2 ~~ 1*C2
P1~~c(e2,e2)*P1
P2~~c(e2,e2)*P2
# covariances
A1 ~~ c(1,.5)*A2
A1 ~~ 0*C1 + 0*C2
A2 ~~ 0*C1 + 0*C2
C1 ~~ c(1,1)*C2"
```

Lets look at some of the critical lines of code in the model:

A1=~ NA*P1 + c(a,a)*P1 Here we create the latent variable A1, the phenotype P for twin 1 (P1) loads on this variable, and in both groups (groups being MZ and DZ twins) the influence of this latent variable on the outcome is the same (contained using `c(a,a)`). Similar code is used to define the latent variables C1 and C2. Now the effect of genes on an outcome is assumed the same for everyone regardless of whether they

are twins, or not, the resemblance between twin 1 and twin 2 difference for MZ and DZ twins. We define/fix the resemblance later in the model here: $A1 \sim c(1, .5) * A2$, because $A1$ and $A2$ are variance 1: $A1 \sim 1 * A1$ and $A2 \sim 1 * A2$ the constrained implies a correlation of 1 for the MZ twins and a correlation of 0.5 for the DZ twins. The common environment is correlated 1 regardless of twin status: $C1 \sim c(1, 1) * C2$, while the unshared environment E is conceptualized as a residual variance of the trait P ($P1$ or $P2$ respectively): $P1 \sim c(e2, e2) * P1$

We assume the latent variables $A(1/2)$ and $C(1/2)$ are uncorrelated, and fix their covariance to 0:

```
A1 ~~ 0*C1 + 0*C2
A2 ~~ 0*C1 + 0*C2
```

We also assume the residual variance (E) is uncorrelated to A and C , but fortunately for us this is a lavaan default. We proceed to fit the model to the simulated data:

```
# Standard ace model:
ace.fit<-cfa(ace.model, data = dataset, group = "zyg")
summary(ace.fit)

## lavaan 0.6-7 ended normally after 20 iterations
##
##      Estimator                      ML
##      Optimization method          NLMINB
##      Number of free parameters      16
##      Number of equality constraints    9
##
##      Number of observations per group:
##      MZ                             1000
##      DZ                             1500
##
## Model Test User Model:
##
##      Test statistic                  3.190
##      Degrees of freedom                3
##      P-value (Chi-square)             0.363
##      Test statistic for each group:
##      MZ                               2.213
##      DZ                               0.977
##
## Parameter Estimates:
##
##      Standard errors                Standard
##      Information                    Expected
##      Information saturated (h1) model Structured
##
## Group 1 [MZ]:
##
## Latent Variables:
##      Estimate Std.Err z-value P(>|z|)
##      A1 =~
##      P1      (a)   0.964   0.075  12.837   0.000
##      A2 =~
##      P2      (a)   0.964   0.075  12.837   0.000
```

```

## C1 =~
## P1      (c)    0.996    0.063    15.755    0.000
## C2 =~
## P2      (c)    0.996    0.063    15.755    0.000
##
## Covariances:
##           Estimate Std.Err z-value P(>|z|)
## A1 ~~
## A2      1.000
## C1      0.000
## C2      0.000
## A2 ~~
## C1      0.000
## C2      0.000
## C1 ~~
## C2      1.000
##
## Intercepts:
##           Estimate Std.Err z-value P(>|z|)
## .P1      0.006    0.055    0.106    0.916
## .P2     -0.025    0.055   -0.461    0.645
## A1      0.000
## A2      0.000
## C1      0.000
## C2      0.000
##
## Variances:
##           Estimate Std.Err z-value P(>|z|)
## A1      1.000
## A2      1.000
## C1      1.000
## C2      1.000
## .P1     (e2)    1.094    0.048    22.698    0.000
## .P2     (e2)    1.094    0.048    22.698    0.000
##
##
## Group 2 [DZ]:
##
## Latent Variables:
##           Estimate Std.Err z-value P(>|z|)
## A1 =~
## P1      (a)    0.964    0.075    12.837    0.000
## A2 =~
## P2      (a)    0.964    0.075    12.837    0.000
## C1 =~
## P1      (c)    0.996    0.063    15.755    0.000
## C2 =~
## P2      (c)    0.996    0.063    15.755    0.000
##
## Covariances:
##           Estimate Std.Err z-value P(>|z|)
## A1 ~~
## A2      0.500
## C1      0.000

```



```
##      C2                0.000
##    A2 ~~
##      C1                0.000
##      C2                0.000
##    C1 ~~
##      C2                1.000
##
## Intercepts:
##              Estimate Std.Err z-value P(>|z|)
##      .P1          -0.077   0.045  -1.724   0.085
##      .P2           0.014   0.045   0.323   0.747
##      A1            0.000
##      A2            0.000
##      C1            0.000
##      C2            0.000
##
## Variances:
##              Estimate Std.Err z-value P(>|z|)
##      A1            1.000
##      A2            1.000
##      C1            1.000
##      C2            1.000
##      .P1      (e2)   1.094   0.048  22.698   0.000
##      .P2      (e2)   1.094   0.048  22.698   0.000
```

And finally lets have a look at the model in a path diagram with tidySEM:

```
lay <- get_layout("", "", "", "", "", "", "",
                  "A1", "", "C1", "", "A2", "", "C2",
                  "", "P1", "", "", "", "P2", "", rows = 3)

graph_sem(model = ace.fit, layout=lay, variance_diameter=.3, angle = 180, rect_height=.35, ellipses_height=.3)
```



```

Ddz <- matrix(c(1,.25,.25,1),2,2) # non-additive genetic correlation for DZ's = 0.5

# make 1000 pairs of MZ twins
MZ <- mvrnorm(1000,mu=c(0,0),Sigma = A+D+E)

# Add a column to label as MZ:
MZ<- cbind.data.frame("MZ",MZ)
colnames(MZ) <- c("zyg","P1", "P2")

# make 1500 DZ twin pairs
DZ <- mvrnorm(1500,mu=c(0,0),Sigma = Adz+Ddz+E)

# add variable too label as DZ:
DZ <- cbind.data.frame("DZ",DZ)
colnames(DZ) <- c("zyg","P1", "P2")

# Combine MZ and DZ twins
dataset <- rbind(MZ,DZ)

```

We then define the lavaan model that can express the variance in the trait P explained by latent variables A, C and E:

```

ade.model<-"
A1=~ NA*P1 + c(a,a)*P1
A2=~ NA*P2 + c(a,a)*P2
D1 =~ NA*P1 + c(d,d)*P1
D2 =~ NA*P2 + c(d,d)*P2
# variances
A1 ~~ 1*A1
A2 ~~ 1*A2
D1 ~~ 1*D1
D2 ~~ 1*D2
P1~~c(e2,e2)*P1
P2~~c(e2,e2)*P2
# covariances
A1 ~~ c(1,.5)*A2
A1 ~~ 0*D1 + 0*D2
A2 ~~ 0*D1 + 0*D2
D1 ~~ c(1,.25)*D2"

```

Lets look at some of the critical lines of code in the model:

A1=~ NA*P1 + c(a,a)*P1 Here we create the latent variable A1, the phenotype P for twin 1 (P1) loads on this variable, and in both groups (groups being MZ and DZ twins) the influence of this latent variable on the outcome is the same (contained using `c(a,a)`). Similar code is used to define the latent variables D1 and D2. Now the effect of genes on an outcome is assumed the same for everyone regardless of whether they are twins, or not, the resemblance between twin 1 and twin 2 difference for MZ and DZ twins. We define/fix the resemblance later in the model here: `A1 ~~ c(1,.5)*A2`, because A1 and A2 are variance 1: `A1 ~~ 1*A1` and `A2 ~~ 1*A2` the constrained implies a correlation of 1 for the MZ twins and a correlation of 0.5 for the DZ twins. The non-additive genetic effects are correlated 1 for MZ twins and .25 for DZ twins `C1 ~~ c(1,.25)*C2`, while the unshared environment E is conceptualized as a residual variance of the trait P (P1 or P2 respectively): `P1~~c(e2,e2)*P1`

We assume the latent variables A(1/2) and D(1/2) are uncorrelated, and fix their covariance to 0:

```
A1 ~~ 0*D1 + 0*D2
A2 ~~ 0*D1 + 0*D2
```

We also assume the residual variance (E) is uncorrelated to A and D, but fortunately for us this is a lavaan default. We proceed to fit the model to the simulated data:

```
# Standard ace model:
ade.fit<-cfa(ade.model, data = dataset, group = "zyg")
summary(ade.fit)
```

```
## lavaan 0.6-7 ended normally after 21 iterations
##
##      Estimator                      ML
##      Optimization method          NLMINB
##      Number of free parameters      16
##      Number of equality constraints    9
##
##      Number of observations per group:
##      MZ                             1000
##      DZ                             1500
##
## Model Test User Model:
##
##      Test statistic                2.717
##      Degrees of freedom              3
##      P-value (Chi-square)           0.437
##      Test statistic for each group:
##      MZ                             1.522
##      DZ                             1.194
##
## Parameter Estimates:
##
##      Standard errors                Standard
##      Information                    Expected
##      Information saturated (h1) model Structured
##
## Group 1 [MZ]:
##
## Latent Variables:
##      Estimate Std.Err z-value P(>|z|)
##      A1 =~
##      P1      (a)   0.833   0.173   4.802   0.000
##      A2 =~
##      P2      (a)   0.833   0.173   4.802   0.000
##      D1 =~
##      P1      (d)   1.101   0.134   8.218   0.000
##      D2 =~
##      P2      (d)   1.101   0.134   8.218   0.000
##
## Covariances:
##      Estimate Std.Err z-value P(>|z|)
##      A1 =~
```

```

##      A2              1.000
##      D1              0.000
##      D2              0.000
##  A2 ~~
##      D1              0.000
##      D2              0.000
##  D1 ~~
##      D2              1.000
##
## Intercepts:
##              Estimate Std.Err  z-value  P(>|z|)
##      .P1           -0.045   0.054   -0.839   0.402
##      .P2            0.005   0.054    0.087   0.931
##      A1            0.000
##      A2            0.000
##      D1            0.000
##      D2            0.000
##
## Variances:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1            1.000
##      A2            1.000
##      D1            1.000
##      D2            1.000
##      .P1      (e2)   1.026   0.045   22.690   0.000
##      .P2      (e2)   1.026   0.045   22.690   0.000
##
##
## Group 2 [DZ]:
##
## Latent Variables:
##              Estimate Std.Err  z-value  P(>|z|)
##  A1 =~
##      P1      (a)    0.833   0.173    4.802   0.000
##  A2 =~
##      P2      (a)    0.833   0.173    4.802   0.000
##  D1 =~
##      P1      (d)    1.101   0.134    8.218   0.000
##  D2 =~
##      P2      (d)    1.101   0.134    8.218   0.000
##
## Covariances:
##              Estimate Std.Err  z-value  P(>|z|)
##  A1 ~~
##      A2              0.500
##      D1              0.000
##      D2              0.000
##  A2 ~~
##      D1              0.000
##      D2              0.000
##  D1 ~~
##      D2              0.250
##
## Intercepts:

```

##		Estimate	Std.Err	z-value	P(> z)	
##	.P1	0.028	0.044	0.627	0.531	
##	.P2	0.060	0.044	1.358	0.175	
##	A1	0.000				
##	A2	0.000				
##	D1	0.000				
##	D2	0.000				
##						
##	Variances:					
##		Estimate	Std.Err	z-value	P(> z)	
##	A1	1.000				
##	A2	1.000				
##	D1	1.000				
##	D2	1.000				
##	.P1	(e2)	1.026	0.045	22.690	0.000
##	.P2	(e2)	1.026	0.045	22.690	0.000

4.3 Sibling interactions

Sibling interactions, or one sibling influencing the others outcome, are an additional mechanism by which twins and siblings can become more alike, or if it concerns a negative sibling interaction become less alike. Sibling interaction cannot be distinguished from rater contrast effects where one child's trait changes to norm or view the rater (usually parent) has of the other child. If one of my children is very quite, the other might seem louder especially in contrast to the other. It is reasonable to assume that sibling interacting effects that persist across self, parental and teacher ratings (especially in the case of twins rated by different teachers) and external or formal measurements are more likely to reflect actual sibling interaction then rater contrast effects.

The "ACE" sibling interaction model (ACE_x) operationalized a phenotype as a function of additive genetic variance (A), common environmental influences(C) (can be rearing, can be societal influences can be governmental policies), and environment unique to each of the individual twins (E) (can be private friends, being in separate classrooms, but also measurement error). The correlation between the genetic influences and common environmental is 1 for MZ twins while the correlation between the genetic influences is 0.5 for DZ twins in addition to these influences the sibling phenotypes are regressed on each other concurrently and the regression in each direction and across MZ and DZ twins is set to be equal.

The ACE sibling interaction model assumes among other things:

1. The absence of non-additive genetic effects (gene gene interaction and dominance)
2. The absence of gene by environment correlation
3. The absence of gene by environment interaction
4. The environments MZ twin grow up are as similar as the environments DZ twins grow up in.

There is also a ADE_x model which is a model one can consider when there is evidence for D in the basic twin model, it is pasted below the example for the ACE_x. The variable "beta" represents the magnitude of the sib interaction, its set high (0.4) and we cranked up the simulated sample size as an ACE_x is very power hungry. In practice you could consider comparing the AEx (settn C to 0) with an ACE model without sibling interaction (setting the interaction to 0) or, if you suspect negative interactions compare the the respective AEx and ADE models.

```
A <- matrix(1,2,2)
C <- matrix(1,2,2)
E <- diag(2)
```

```

Adz <- matrix(c(1,.5,.5,1),2,2)

#sibling interactie effect:
beta <- .4

# MZ twins

MZ <- mvrnorm(4000,mu=c(0,0),Sigma = A+C+E)

# regress the sibs on eachother
MZ <- t(matrix(c(1,beta,beta,1),2,2) %*% t(MZ))

# MZ data frame mwithcolumns zygotity ("MZ"), twin1 data, twin2 data
MZ <- cbind.data.frame("MZ",MZ)
colnames(MZ) <- c("zyg", "P1", "P2")

# DZ twins
DZ <- mvrnorm(6000,mu=c(0,0),Sigma = Adz+C+E)

# add sib interaction
DZ <- t(matrix(c(1,beta ,beta ,1),2,2) %*% t(DZ))

DZ <- cbind.data.frame("DZ",DZ)

colnames(DZ) <- c("zyg", "P1", "P2")

# cmbineer MZ en DZ in een dataset:
dataset <- rbind(MZ,DZ)

```

Having generated data with a sibling interaction we can go ahead and fit the sibling interaction model below to retrieve the simulate parameters.

```

# Models

# Sibling interaction model:
ace.model.sib.int<-"
A1=~ NA*P1 + c(a,a)*P1
A2=~ NA*P2 + c(a,a)*P2
C1 =~ NA*P1 + c(c,c)*P1
C2 =~ NA*P2 + c(c,c)*P2
# variances
A1 ~~ 1*A1
A2 ~~ 1*A2
C1 ~~ 1*C1
C2 ~~ 1*C2
P1~~c(e2,e2)*P1
P2~~c(e2,e2)*P2
# covariances
A1 ~~ c(1,.5)*A2
A1 ~~ 0*C1 + 0*C2
A2 ~~ 0*C1 + 0*C2
C1 ~~ c(1,1)*C2

```

```

# regs
P1 ~ c(beta,beta)*P2
P2 ~ c(beta,beta)*P1
"

# Sibling interaction:
ace.fitsib <-cfa(ace.model.sib.int, data = dataset,group = "zyg")
summary(ace.fitsib)

```

```

## lavaan 0.6-7 ended normally after 44 iterations
##
##      Estimator                      ML
##      Optimization method          NLMINB
##      Number of free parameters      20
##      Number of equality constraints   12
##
##      Number of observations per group:
##      MZ                             4000
##      DZ                             6000
##
## Model Test User Model:
##
##      Test statistic                  0.625
##      Degrees of freedom                2
##      P-value (Chi-square)             0.732
##      Test statistic for each group:
##      MZ                             0.208
##      DZ                             0.417
##
## Parameter Estimates:
##
##      Standard errors                Standard
##      Information                    Expected
##      Information saturated (h1) model Structured
##
## Group 1 [MZ]:
##
## Latent Variables:
##      Estimate Std.Err z-value P(>|z|)
##      A1 =~
##      P1      (a)   0.858   0.033  25.688   0.000
##      A2 =~
##      P2      (a)   0.858   0.033  25.688   0.000
##      C1 =~
##      P1      (c)   0.509   0.251   2.029   0.042
##      C2 =~
##      P2      (c)   0.509   0.251   2.029   0.042
##
## Regressions:
##      Estimate Std.Err z-value P(>|z|)
##      P1 ~
##      P2      (beta)  0.475   0.043  11.099   0.000

```



```

## P2 ~
## P1 (beta) 0.475 0.043 11.099 0.000
##
## Covariances:
## Estimate Std.Err z-value P(>|z|)
## A1 ~~
## A2 1.000
## C1 0.000
## C2 0.000
## A2 ~~
## C1 0.000
## C2 0.000
## C1 ~~
## C2 1.000
##
## Intercepts:
## Estimate Std.Err z-value P(>|z|)
## .P1 0.003 0.021 0.139 0.889
## .P2 -0.007 0.021 -0.350 0.726
## A1 0.000
## A2 0.000
## C1 0.000
## C2 0.000
##
## Variances:
## Estimate Std.Err z-value P(>|z|)
## A1 1.000
## A2 1.000
## C1 1.000
## C2 1.000
## .P1 (e2) 0.797 0.053 15.019 0.000
## .P2 (e2) 0.797 0.053 15.019 0.000
##
##
## Group 2 [DZ]:
##
## Latent Variables:
## Estimate Std.Err z-value P(>|z|)
## A1 =~
## P1 (a) 0.858 0.033 25.688 0.000
## A2 =~
## P2 (a) 0.858 0.033 25.688 0.000
## C1 =~
## P1 (c) 0.509 0.251 2.029 0.042
## C2 =~
## P2 (c) 0.509 0.251 2.029 0.042
##
## Regressions:
## Estimate Std.Err z-value P(>|z|)
## P1 ~
## P2 (beta) 0.475 0.043 11.099 0.000
## P2 ~
## P1 (beta) 0.475 0.043 11.099 0.000
##

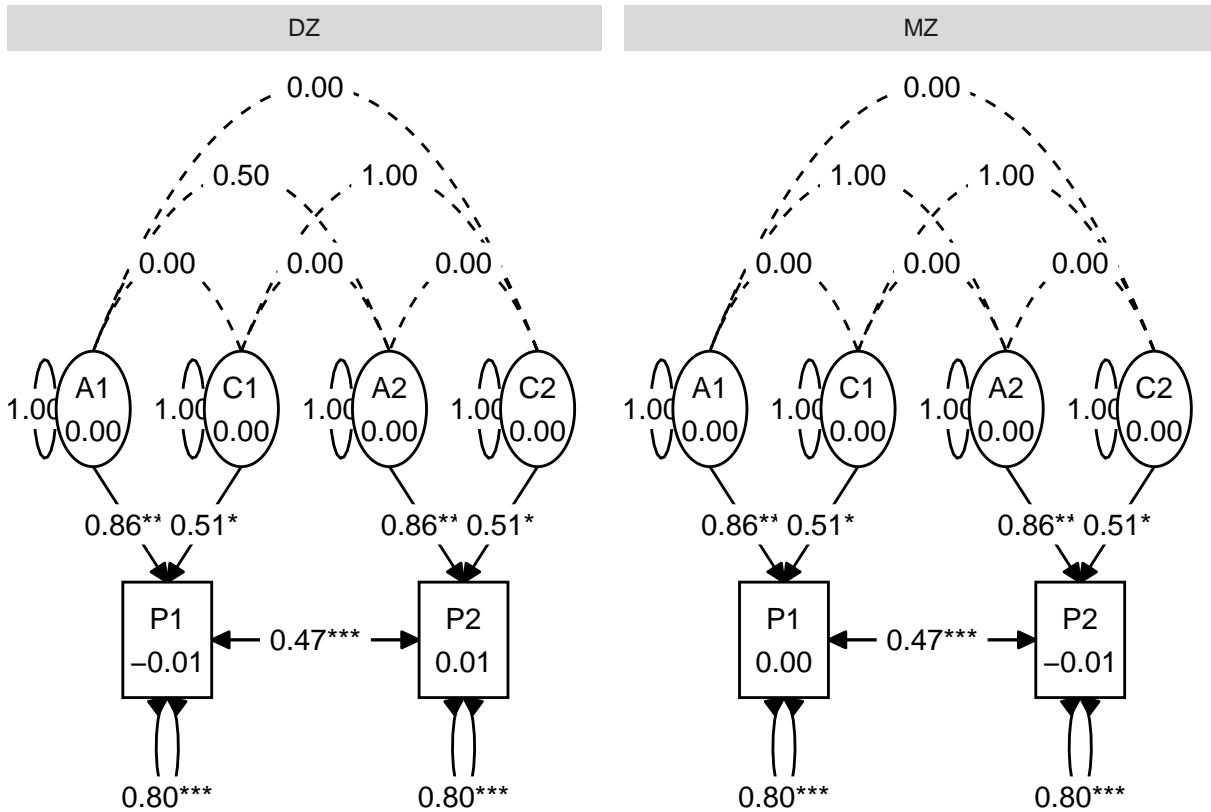
```

```
## Covariances:
##           Estimate Std.Err z-value P(>|z|)
##   A1 ~~
##     A2           0.500
##     C1           0.000
##     C2           0.000
##   A2 ~~
##     C1           0.000
##     C2           0.000
##   C1 ~~
##     C2           1.000
##
## Intercepts:
##           Estimate Std.Err z-value P(>|z|)
##     .P1          -0.010   0.017  -0.582   0.560
##     .P2           0.008   0.017   0.477   0.633
##     A1           0.000
##     A2           0.000
##     C1           0.000
##     C2           0.000
##
## Variances:
##           Estimate Std.Err z-value P(>|z|)
##     A1           1.000
##     A2           1.000
##     C1           1.000
##     C2           1.000
##     .P1 (e2)      0.797   0.053  15.019   0.000
##     .P2 (e2)      0.797   0.053  15.019   0.000
```

Lets proceed to visualize the path diagram

```
lay <- get_layout("A1","", "C1","", "A2","", "C2",
                  "", "P1","", "", "P2","", rows = 2)

graph_sem(model = ace.fitsib, layout=lay, variance_diameter=.3, angle = 180, rect_height=.35, ellipses_height=.35)
```



And for completeness here is the ADE sibling interaction model:

```
# Sibling interaction model:
ade.model.sib.int<-"
A1=~ NA*P1 + c(a,a)*P1
A2=~ NA*P2 + c(a,a)*P2
D1 =~ NA*P1 + c(c,c)*P1
D2 =~ NA*P2 + c(c,c)*P2
# variances
A1 ~~ 1*A1
A2 ~~ 1*A2
D1 ~~ 1*D1
D2 ~~ 1*D2
P1~~c(e2,e2)*P1
P2~~c(e2,e2)*P2
# covariances
A1 ~~ c(1,.5)*A2
A1 ~~ 0*D1 + 0*D2
A2 ~~ 0*D1 + 0*D2
D1 ~~ c(1,.25)*D2

# regs
P1 ~ c(beta,beta)*P2
P2 ~ c(beta,beta)*P1
"
```

4.4 Binary/Ordinal data

In lavaan binary or ordinal data can be readily analyzed by specifying the variables are ordered, Lets simulate twin data generated under and ACE model, from multivariate normal data that is then transformed to ordered data.

recall the ACE model assumes among other things:

1. The absence of non-additive genetic effects (gene gene interaction and dominance)
2. The absence of sibling influences or rater contrast effect
3. The absence of gene by environment correlation
4. The absence of gene by environment interaction
5. The environments MZ twin grow up are as similar the environments DZ twins grow up in.

The binary data model assume the following:

1. A normally distributed continuous latent variable causes the observed ordinal (or binary) variable.

4.4.1 Ordinal data

Lets simulate same data where the resemblance between twins is a function of equal parts A,C and E (33.3%/33.3%/33.4%) and the data is ordered in nature:

```
A <- matrix(1,2,2) # genetic correlation for MZ's = 1
C <- matrix(1,2,2)
E <- diag(2)
Adz <- matrix(c(1,.5,.5,1),2,2) # genetic correlation for DZ's = 0.5

# make 1000 pairs of MZ twins
MZ <- mvrnorm(1000,mu=c(0,0),Sigma = A+C+E)

# Add a column to label as MZ:
MZ<- cbind.data.frame("MZ",MZ)
colnames(MZ) <- c("zyg","P1", "P2")

# make 1500 DZ twin pairs
DZ <- mvrnorm(1500,mu=c(0,0),Sigma = Adz+C+E)

# add variable too label as DZ:
DZ <- cbind.data.frame("DZ",DZ)
colnames(DZ) <- c("zyg","P1", "P2")

# Combine MZ and DZ twins
dataset <- rbind(MZ,DZ)

# make the data ordered:
dataset[dataset[,2] < 0 ,2] <- 0
dataset[dataset[,2] > 0 & dataset[,2] < 1,2] <- 1
dataset[dataset[,2] > 1 ,2] <- 2

dataset[dataset[,3] < 0 ,3] <- 0
dataset[dataset[,3] > 0 & dataset[,3] < 1 ,3] <- 1
dataset[dataset[,3] > 1 ,3] <- 2
```

We then define the lavaan model that can express the variance in the trait P explained by latent variables A, C and E:

```
ace.model<-"
A1=~ NA*P1 + c(a,a)*P1
A2=~ NA*P2 + c(a,a)*P2
C1 =~ NA*P1 + c(c,c)*P1
C2 =~ NA*P2 + c(c,c)*P2
# variances
A1 ~~ 1*A1
A2 ~~ 1*A2
C1 ~~ 1*C1
C2 ~~ 1*C2
P1~~c(e2,e2)*P1
P2~~c(e2,e2)*P2
# covariances
A1 ~~ c(1,.5)*A2
A1 ~~ 0*C1 + 0*C2
A2 ~~ 0*C1 + 0*C2
C1 ~~ c(1,1)*C2

# first threshold fixed:
P1 | 0*t1 + c(t,t)*t2
P2 | 0*t1 + c(t,t)*t2
"
```

Lets look at some of the critical lines of code in the model:

IF we analyze ordinal data in a SEM model we model a latent (normaly distributed) variable that is the “cause” of the observed ordinal variable. This latent variable has various features, it has a mean, a variance and there ar thresholds, which are the values of the latent continuous variable at which the observed variable increases from 0 to 1 (threshold 1) from 1 to 2 (threshold 1) and further if we have more ordered categories.

To identify a model with ordinal variables, we have to chose to either estimate the thresholds, or the mean and variance. From a twin modeling perspective, it makes more sense to estimate the variance, as we wish to partition the variance, therefore the code below fixes the first thresholds in the data to 0, we have no need or wish to estimate the means off the latent variables (assumed 0) so we need only fix 1 threshold.:

```
# first thresholds fixed:
P1 | 0*t1 + c(t,t)*t2
P2 | 0*t1 + c(t,t)*t2
```

Then we can proceed to fit the model, note we add the following arguments to the `cfa()` function to let lavaan know the data is ordered, and we wish to use a specific “parameterization” the details of which are beyond the scope of the current document, but it is essential to specify this argument `cfa(... ,parameterization="theta",ordered=TRUE)`.

```
# Standard ace model:
ace.fit<-cfa(ace.model, data = dataset,group = "zyg",parameterization="theta",ordered=TRUE)
```

```
## Warning in lav_model_vcov(lavmodel = lavmodel, lavsamplestats = lavsamplestats, : lavaan WARNING:
## The variance-covariance matrix of the estimated parameters (vcov)
## does not appear to be positive definite! The smallest eigenvalue
## (= 7.803844e-16) is close to zero. This may be a symptom that the
## model is not identified.
```

```
summary(ace.fit)
```

```
## lavaan 0.6-7 ended normally after 21 iterations
##
## Estimator DWLS
## Optimization method NLMINB
## Number of free parameters 16
## Number of equality constraints 12
##
## Number of observations per group:
## MZ 1000
## DZ 1500
##
## Model Test User Model:
## Standard Robust
## Test Statistic 6.509 5.929
## Degrees of freedom 6 6
## P-value (Chi-square) 0.369 0.431
## Scaling correction factor 1.216
## Shift parameter for each group:
## MZ 0.231
## DZ 0.346
## simple second-order correction
## Test statistic for each group:
## MZ 5.984 5.151
## DZ 0.524 0.778
##
## Parameter Estimates:
##
## Standard errors Robust.sem
## Information Expected
## Information saturated (h1) model Unstructured
##
## Group 1 [MZ]:
##
## Latent Variables:
## Estimate Std.Err z-value P(>|z|)
## A1 =~
## P1 (a) 1.166 0.130 8.937 0.000
## A2 =~
## P2 (a) 1.166 0.130 8.937 0.000
## C1 =~
## P1 (c) 1.101 0.093 11.831 0.000
## C2 =~
## P2 (c) 1.101 0.093 11.831 0.000
##
## Covariances:
## Estimate Std.Err z-value P(>|z|)
## A1 ~~
## A2 1.000
## C1 0.000
## C2 0.000
```

```

##      A2 ~~
##      C1          0.000
##      C2          0.000
##      C1 ~~
##      C2          1.000
##
## Intercepts:
##              Estimate Std.Err  z-value  P(>|z|)
##      .P1          0.000
##      .P2          0.000
##      A1          0.000
##      A2          0.000
##      C1          0.000
##      C2          0.000
##
## Thresholds:
##              Estimate Std.Err  z-value  P(>|z|)
##      P1|t1          0.000
##      P1|t2      (t)    1.100    0.039   27.940    0.000
##      P2|t1          0.000
##      P2|t2      (t)    1.100    0.039   27.940    0.000
##
## Variances:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1          1.000
##      A2          1.000
##      C1          1.000
##      C2          1.000
##      .P1      (e2)    1.070    0.053   20.301    0.000
##      .P2      (e2)    1.070    0.053   20.301    0.000
##
## Scales y*:
##              Estimate Std.Err  z-value  P(>|z|)
##      P1          0.524
##      P2          0.524
##
## Group 2 [DZ]:
##
## Latent Variables:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1 =~
##      P1      (a)    1.166    0.130    8.937    0.000
##      A2 =~
##      P2      (a)    1.166    0.130    8.937    0.000
##      C1 =~
##      P1      (c)    1.101    0.093   11.831    0.000
##      C2 =~
##      P2      (c)    1.101    0.093   11.831    0.000
##
## Covariances:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1 =~
##      A2          0.500

```

```

##      C1              0.000
##      C2              0.000
##    A2 ~~
##      C1              0.000
##      C2              0.000
##    C1 ~~
##      C2              1.000
##
## Intercepts:
##              Estimate Std.Err  z-value  P(>|z|)
##      .P1              0.000
##      .P2              0.000
##      A1              0.000
##      A2              0.000
##      C1              0.000
##      C2              0.000
##
## Thresholds:
##              Estimate Std.Err  z-value  P(>|z|)
##      P1|t1              0.000
##      P1|t2      (t)    1.100    0.039   27.940    0.000
##      P2|t1              0.000
##      P2|t2      (t)    1.100    0.039   27.940    0.000
##
## Variances:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1              1.000
##      A2              1.000
##      C1              1.000
##      C2              1.000
##      .P1      (e2)    1.070    0.053   20.301    0.000
##      .P2      (e2)    1.070    0.053   20.301    0.000
##
## Scales y*:
##              Estimate Std.Err  z-value  P(>|z|)
##      P1              0.524
##      P2              0.524

```

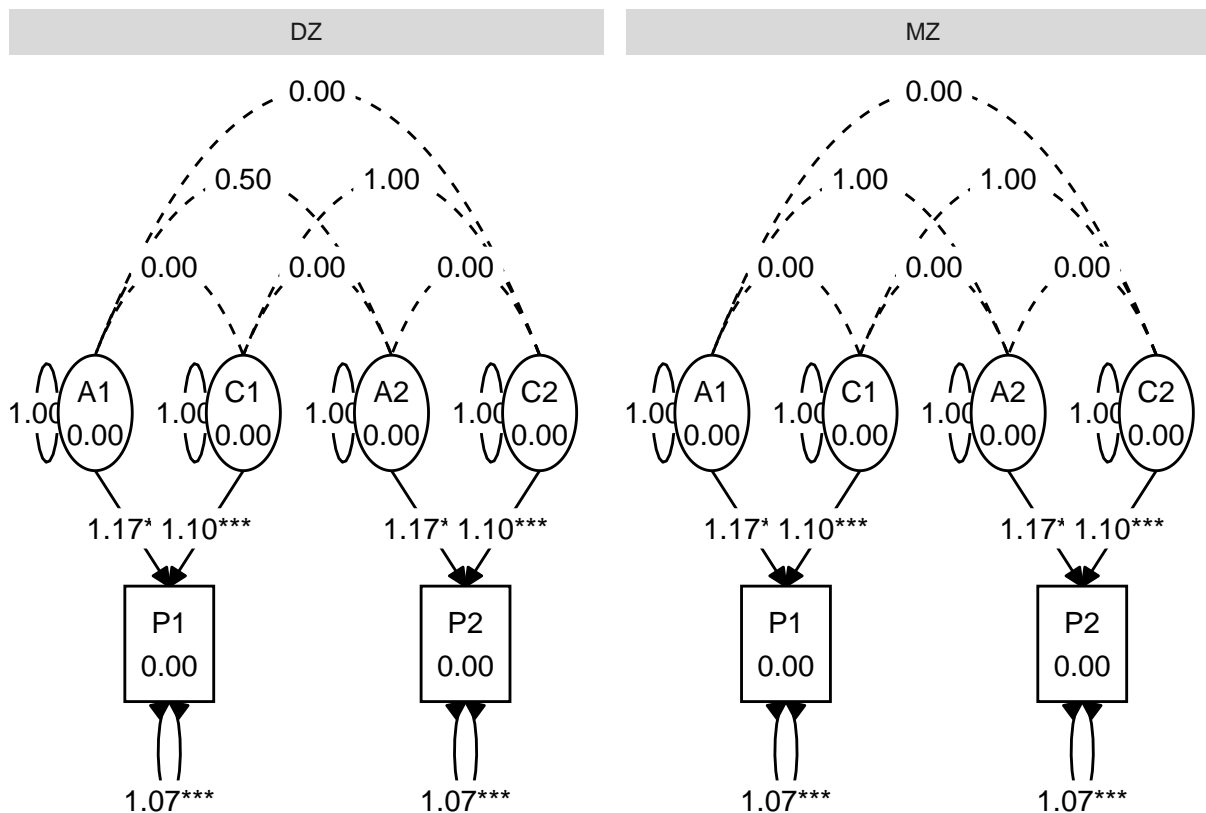
And finally lets have a look at the model in a path diagram with tidySEM:

```

lay <- get_layout("", "", "", "", "", "", "",
                  "A1", "", "C1", "", "A2", "", "C2",
                  "", "P1", "", "", "", "P2", "", rows = 3)

graph_sem(model = ace.fit, layout=lay, variance_diameter=.3, angle = 180, rect_height=.35, ellipses_height=.3)

```

4.4.2 Binary data

Lets simulate same data where the resemblance between twins is a function of equal parts A,C and E (33.3%/33.3%/33.4%) and the data is binary in nature (so case/control for example):

```
A <- matrix(1,2,2) # genetic correlation for MZ's = 1
C <- matrix(1,2,2)
E <- diag(2)
Adz <- matrix(c(1,.5,.5,1),2,2) # genetic correlation for DZ's = 0.5

# make 1000 pairs of MZ twins
MZ <- mvrnorm(1000,mu=c(0,0),Sigma = A+C+E)

# Add a column to label as MZ:
MZ<- cbind.data.frame("MZ",MZ)
colnames(MZ) <- c("zyg","P1", "P2")

# make 1500 DZ twin pairs
DZ <- mvrnorm(1500,mu=c(0,0),Sigma = Adz+C+E)

# add variable too label as DZ:
DZ <- cbind.data.frame("DZ",DZ)
colnames(DZ) <- c("zyg","P1", "P2")
```

```

# Combine MZ and DZ twins
dataset <- rbind(MZ,DZ)

# make the data ordered:
dataset[dataset[,2] < 1 ,2] <- 0
dataset[dataset[,2] > 1,2] <- 1

dataset[dataset[,3] < 1 ,3] <- 0
dataset[dataset[,3] > 1 ,3] <- 1

```

We then define the lavaan model that can express the variance in the trait P explained by latent variables A, C and E:

```

ace.model<-"
A1=~ NA*P1 + c(a,a)*P1
A2=~ NA*P2 + c(a,a)*P2
C1 =~ NA*P1 + c(c,c)*P1
C2 =~ NA*P2 + c(c,c)*P2
# variances
A1 ~~ 1*A1
A2 ~~ 1*A2
C1 ~~ 1*C1
C2 ~~ 1*C2
P1~~c(e2,e2)*P1
P2~~c(e2,e2)*P2
# covariances
A1 ~~ c(1,.5)*A2
A1 ~~ 0*C1 + 0*C2
A2 ~~ 0*C1 + 0*C2
C1 ~~ c(1,1)*C2

# threshold fixed:
P1 | 1*t1
P2 | 1*t1
"

```

Lets look at some of the critical lines of code in the model:

IF we analyze ordinal data in a SEM model we model a latent (normaly distributed) variable that is the “cause” of the observed ordinal variable. This latent variable has various features, it has a mean, a variance and there ar thresholds, which are the values of the latent continuous variable at which the observed variable increases from 0 to 1 (threshold 1) from 1 to 2 (threshold 1) and further if we have more ordered categories.

To identify a model with binary outcome, we have to chose to either estimate the threshold, or the variance. From a twin modeling perspective, it makes more sense to estimate the variance, as we wish to partition the variance, therefore the code below fixes the threshold to 0:

```

# thresholds fixed:
P1 | 1*t1
P2 | 1*t1

```

Then we can proceed to fit the model, note we add the following arguments to the `cfa()` function to let lavaan know the data is ordered, and we wish to use a specific “parameterization” the details of which

are beyond the scope of the current document, but it is essential to specify this argument `cfa(... ,parameterization="theta",ordered=TRUE)`.

Standard ace model:

```
ace.fit<-cfa(ace.model, data = dataset,group = "zyg",parameterization="theta",ordered=TRUE)
summary(ace.fit)
```

```
## lavaan 0.6-7 ended normally after 20 iterations
##
##      Estimator                      DWLS
##      Optimization method          NLMINB
##      Number of free parameters      12
##      Number of equality constraints    9
##
##      Number of observations per group:
##      MZ                             1000
##      DZ                             1500
##
## Model Test User Model:
##
##      Standard      Robust
##      Test Statistic    0.645    0.428
##      Degrees of freedom      3      3
##      P-value (Chi-square)    0.886    0.934
##      Scaling correction factor      0.774
##      Shift parameter for each group:
##      MZ                             -0.162
##      DZ                             -0.244
##      simple second-order correction
##      Test statistic for each group:
##      MZ                             0.412    0.370
##      DZ                             0.233    0.058
##
## Parameter Estimates:
##
##      Standard errors          Robust.sem
##      Information              Expected
##      Information saturated (h1) model    Unstructured
##
## Group 1 [MZ]:
##
## Latent Variables:
##      Estimate Std.Err z-value P(>|z|)
##      A1 =~
##      P1      (a)   0.757   0.179   4.217   0.000
##      A2 =~
##      P2      (a)   0.757   0.179   4.217   0.000
##      C1 =~
##      P1      (c)   1.022   0.109   9.386   0.000
##      C2 =~
##      P2      (c)   1.022   0.109   9.386   0.000
##
## Covariances:
##      Estimate Std.Err z-value P(>|z|)
```

```

##      A1 ~~
##      A2          1.000
##      C1          0.000
##      C2          0.000
##      A2 ~~
##      C1          0.000
##      C2          0.000
##      C1 ~~
##      C2          1.000
##
## Intercepts:
##              Estimate Std.Err  z-value  P(>|z|)
##      .P1            0.000
##      .P2            0.000
##      A1            0.000
##      A2            0.000
##      C1            0.000
##      C2            0.000
##
## Thresholds:
##              Estimate Std.Err  z-value  P(>|z|)
##      P1|t1          1.000
##      P2|t1          1.000
##
## Variances:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1            1.000
##      A2            1.000
##      C1            1.000
##      C2            1.000
##      .P1      (e2)  0.919    0.115    8.007    0.000
##      .P2      (e2)  0.919    0.115    8.007    0.000
##
## Scales y*:
##              Estimate Std.Err  z-value  P(>|z|)
##      P1            0.628
##      P2            0.628
##
##
## Group 2 [DZ]:
##
## Latent Variables:
##              Estimate Std.Err  z-value  P(>|z|)
##      A1 =~
##      P1      (a)    0.757    0.179    4.217    0.000
##      A2 =~
##      P2      (a)    0.757    0.179    4.217    0.000
##      C1 =~
##      P1      (c)    1.022    0.109    9.386    0.000
##      C2 =~
##      P2      (c)    1.022    0.109    9.386    0.000
##
## Covariances:
##              Estimate Std.Err  z-value  P(>|z|)

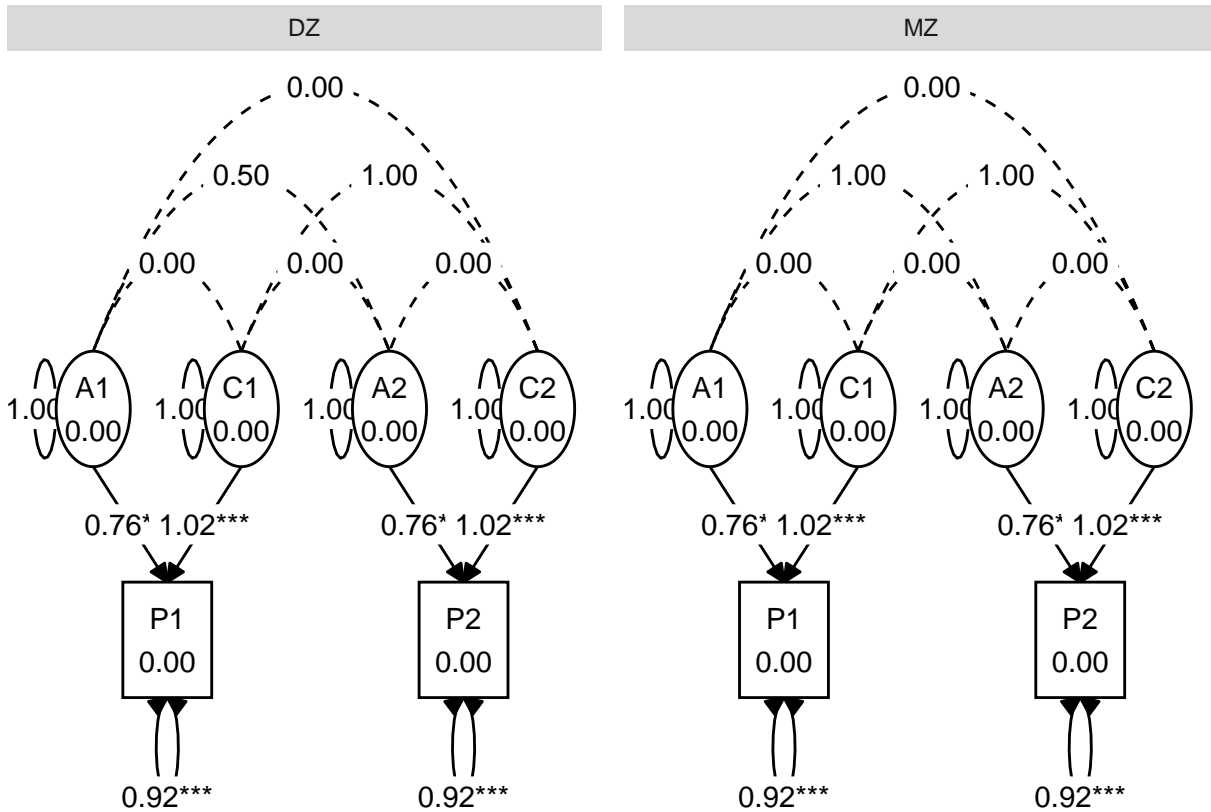
```

```
## A1 ~~
## A2      0.500
## C1      0.000
## C2      0.000
## A2 ~~
## C1      0.000
## C2      0.000
## C1 ~~
## C2      1.000
##
## Intercepts:
##           Estimate Std.Err z-value P(>|z|)
## .P1          0.000
## .P2          0.000
## A1           0.000
## A2           0.000
## C1           0.000
## C2           0.000
##
## Thresholds:
##           Estimate Std.Err z-value P(>|z|)
## P1|t1        1.000
## P2|t1        1.000
##
## Variances:
##           Estimate Std.Err z-value P(>|z|)
## A1           1.000
## A2           1.000
## C1           1.000
## C2           1.000
## .P1 (e2)    0.919   0.115   8.007   0.000
## .P2 (e2)    0.919   0.115   8.007   0.000
##
## Scales y*:
##           Estimate Std.Err z-value P(>|z|)
## P1          0.628
## P2          0.628
```

And finally lets have a look at the model in a path diagram with tidySEM:

```
lay <- get_layout("", "", "", "", "", "", "",
                  "A1", "", "C1", "", "A2", "", "C2",
                  "", "P1", "", "", "", "P2", "", rows = 3)

graph_sem(model = ace.fit, layout=lay, variance_diameter=.3, angle = 180, rect_height=.35, ellipses_height=.3)
```



4.5 Sex specific effects

4.6 Gene-environment interaction

The most intuitive models for gene-environment interaction with a continuous environment(Purcell 2002) cannot be fitted in lavaan, please consider OpenMx or Mplus. The issue is that the model requires a factor loading to be a linear function of an observed variable (moderation) and this isn't implemented in lavaan.

4.7 Gene-environment correlation

4.8 Rater bias models

5 Multitrait twin models in lavaan

5.1 Cholesky decomposition

5.2 Direction of Causation models

6 Longitudinal twin models

6.1 Auto regressive models

6.1.1 Phenotype to environment effects

6.2 Growth curve models

References

- Bates, Timothy C., Hermine Maes, and Michael C. Neale. 2019. “Umx: Twin and Path-Based Structural Equation Modeling in R.” *Twin Research and Human Genetics* 22 (1): 27–41. <https://doi.org/10.1017/thg.2019.2>.
- Huang, Po-Hsien, Hung Chen, and Li-Jen Weng. 2017. “A Penalized Likelihood Method for Structural Equation Modeling.” *Psychometrika* 82 (2): 329–54. <https://doi.org/10.1007/s11336-017-9566-9>.
- Jacobucci, Ross, Kevin J. Grimm, and John J. McArdle. 2016. “Regularized Structural Equation Modeling.” *Structural Equation Modeling: A Multidisciplinary Journal* 23 (4): 555–66. <https://doi.org/10.1080/10705511.2016.1154793>.
- Merkle, Edgar C., and Yves Rosseel. 2018. “Blavaan: Bayesian Structural Equation Models via Parameter Expansion.” *Journal of Statistical Software* 85 (4). <https://doi.org/10.18637/jss.v085.i04>.
- Neale, Michael C., Michael D. Hunter, Joshua N. Pritikin, Mahsa Zahery, Timothy R. Brick, Robert M. Kirkpatrick, Ryne Estabrook, Timothy C. Bates, Hermine H. Maes, and Steven M. Boker. 2015. “OpenMx 2.0: Extended Structural Equation and Statistical Modeling.” *Psychometrika* 81 (2): 535–49. <https://doi.org/10.1007/s11336-014-9435-8>.
- Purcell, Shaun. 2002. “Variance Components Models for GeneEnvironment Interaction in Twin Analysis.” *Twin Research* 5 (6): 554–71. <https://doi.org/10.1375/136905202762342026>.
- Rosseel, Yves. 2012. “Lavaan: AnRPackage for Structural Equation Modeling.” *Journal of Statistical Software* 48 (2). <https://doi.org/10.18637/jss.v048.i02>.