

SEM 1: Confirmatory Factor Analysis

Week 4 - Missing data, EFA and higher order models

Sacha Epskamp

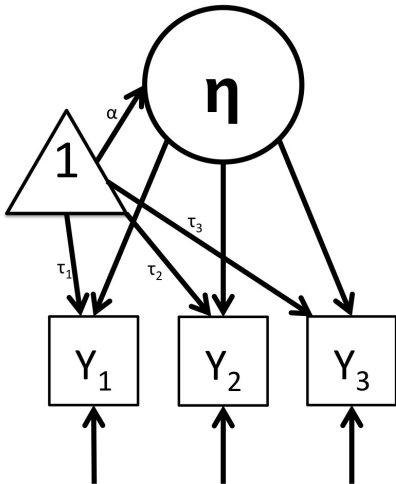
18-04-2017

Mean structure

$$\Sigma = \Lambda\Psi\Lambda^T + \Theta$$

$$\mu = \tau + \Lambda\alpha$$

- τ can cancel α out, hence we need to identify $\alpha = \mathbf{0}$
- Number of parameters: $p(p+1)/2$ variances and covariances and p means!
- Number of parameters: p intercepts in τ
- p more observations, and p more parameters. This is why we normally ignore means!



Steps to assess measurement invariance:

- **Configural invariance:** Is the configuration of the model the same?
- **Weak Invariance:** Are factor loadings the same?
- **Strong Invariance:** Are the intercepts the same?
- **Strict Invariance:** Are the residual variances the same?

I see myself as someone who is talkative

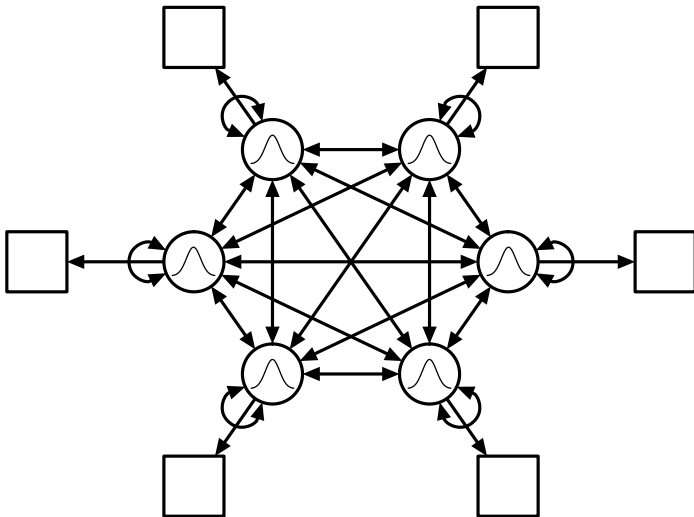
Disagree
strongly
1

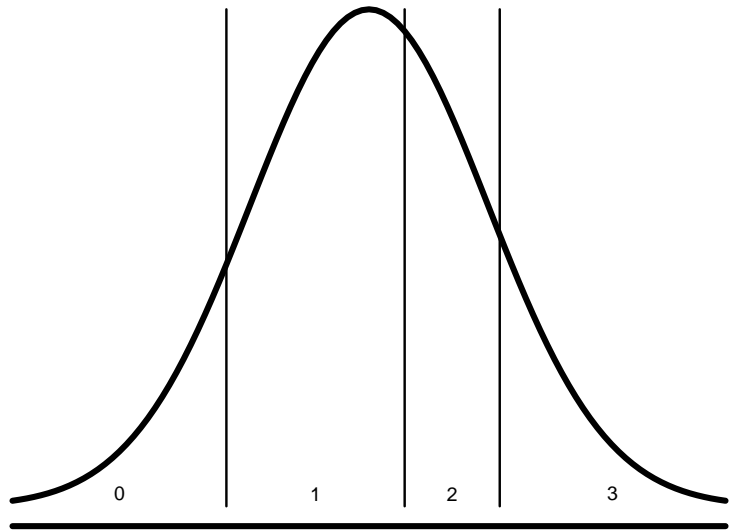
Disagree
a little
2

Neither agree
nor disagree
3

Agree
a little
4

Agree
Strongly
5





Sample Size

How big is 'big enough'?

- $n : q$ ratio should be high
 - Theory: to efficiently estimate lots of parameters, a larger sample is needed (5-10 per parameter)
 - There's very little evidence that it matters (Jackson, 2003)
 - This ratio is less important than absolute sample size
- $n \approx 200$ people
 - This is median SEM sample size (Shah & Goldstein, 2006)
 - Appropriate for an average model with ML estimation
 - Other recommendations: 100-200 people minimum
- Use larger n if:
 - Assumptions are violated (e.g., data are nonnormal)
 - Model is complex (e.g., latent interactions, multilevel structure)
 - Indicators have low reliability (factor loadings are low)

Power for Test of (Not-)Close Fit

- RMSEA estimates a population value
 - Its sampling distribution has been worked out
 - So we can put a confidence interval around it
 - This confidence interval allows us to ask whether RMSEA is significantly different from a specified value
- If the population model fit is NOT CLOSE, what is power to reject H_0 by the test of close fit?
- If the population model fit is CLOSE, what is power to reject H_0 by the test of not-close fit?
- Method described in MacCallum et al. (1996) is implemented in online calculators:
 - Power and minimum sample size for RMSEA:
<http://quantpsy.org/rmsear/rmsear.htm>
 - Power curves for RMSEA:
<http://quantpsy.org/rmsear/rmsearplot.htm>
 - See also `findRMSEAsamplesize` in `semTools`

Why are data missing? In a general X predicts Y case:

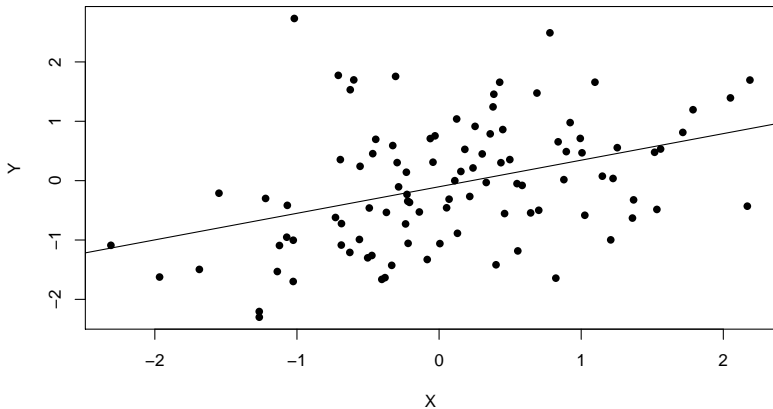
- **Missing completely at random (MCAR)**
 - Missingness is independent of Y or X
 - Everything is fine!
- **Missing at random (MAR)**
 - Missingness is independent of Y , but not of X
 - Example: Men less willing to respond to mental health questionnaire
 - Not a big problem
- **Missing not at random (MNAR)**
 - Missingness depends on Y
 - Example: People with severe mental health problems fill in less questionnaires
 - This is bad :(

Unfortunately,

A dataset:

X	Y
5	5
6	5
5	6
8	5
6	7
7	7
6	9
9	8
9	9
12	9

A larger dataset:

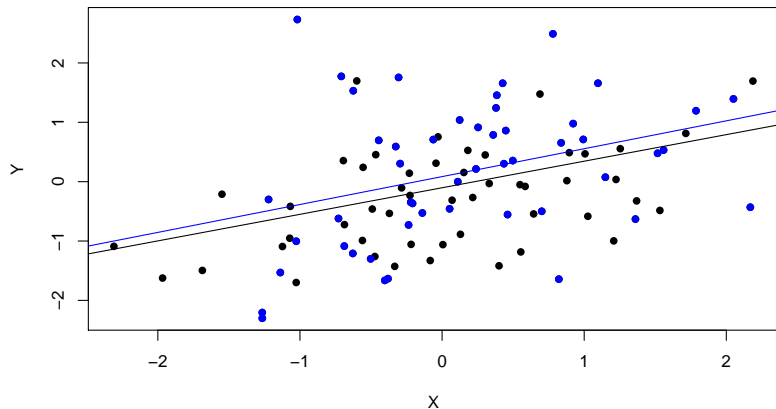


Missing completely at random (MCAR):

X	Y
5	5
6	5
5	6
8	5
6	7
7	7
6	9
9	8
9	9
12	9

MCAR

A larger dataset:



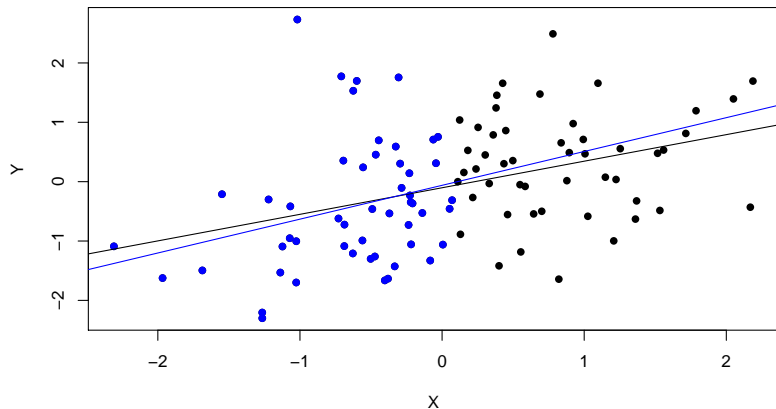
Blue = observed

Missing at random (MAR):

X	Y
5	5
6	5
5	6
8	5
6	7
7	7
6	9
9	8
9	9
12	9

MAR

A larger dataset:

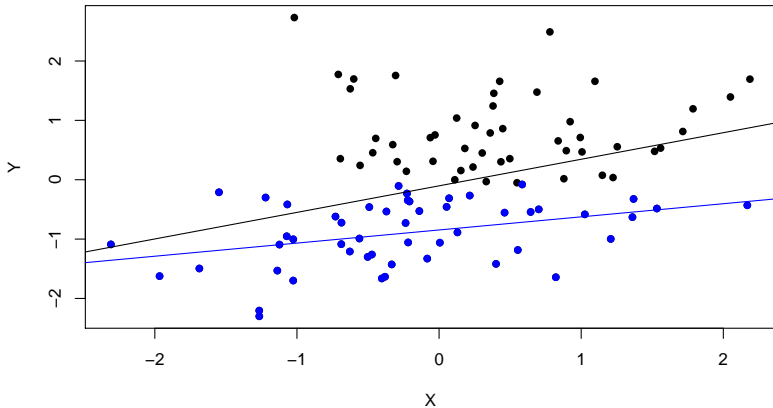


Blue = observed

Missing not at random (MNAR):

X	Y
5	5
6	5
5	6
8	5
6	7
7	7
6	9
9	8
9	9
12	9

MNAR



Blue = observed

Missing data

- Best case: no missings
- MCAR or MAR: this is ok
- MNAR: This is not ok
- Unfortunately, no real statistical way to checking if missings are MNAR
- Thus, MAR needs to be assumed to continue

Old ways of handling missing data

- Compute \mathbf{S} using list-wise deletion
 - Delete all rows with a missing value
 - Downside: deletes observed data
- Compute \mathbf{S} using pair-wise estimation
 - Estimate each element of \mathbf{S} using all available data
 - Downside: Each covariance is based on different n
- (multiple) imputations
 - Impute missingness using mean scores or regression models
 - Downside: complicated, can increase bias if MNAR

Modern way: full-information maximum likelihood (FIML)

- Uses the full data set and all observations
- Downside: full data needed (analysis can not be done using covariance matrix)
- Assumes MAR

```
fit <- cfa(model, data, missing = "FIML")
```

Assumptions of ML

1. Independence: Observations are a simple random sample from some population
 - Consequence: underestimated standard errors, inflated Type-I error rates
 - Solution 1: use SE correction for dependence structure
 - Solution 2: multilevel SEM
2. Multivariate Normality: Variables are univariate normally distributed at levels of all other variables, residuals are normal and homoscedastic, latent variables are normal, bivariate relations are linear
 - Consequence: standard errors are incorrect (probably too low), Type-I error rate is not accurate (probably too high)
 - Solution 1: use robust standard errors (`estimator = 'MLM'`, with complete data; `estimator = 'MLR'` with incomplete data)
 - Solution 2: use bootstrapped standard errors & test statistic

Exploratory Factor Analysis (EFA)

Exploratorily estimate Λ (no free elements in Λ):

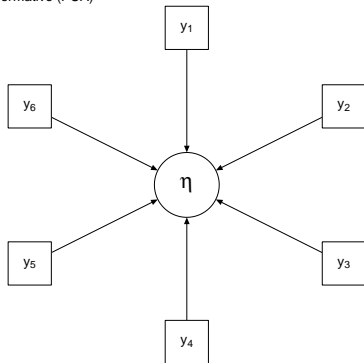
$$\Sigma = \Lambda\Psi\Lambda^T + \Theta$$

Very close, but not the same (!!) as principal component analysis (PCA):

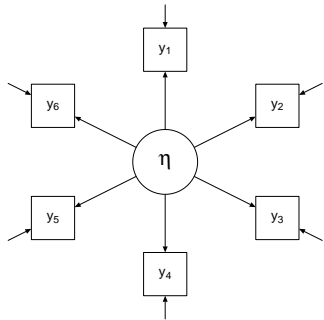
$$\Sigma = \Lambda\Psi\Lambda^T$$

Very different interpretation. EFA *measures* latents (there is measurement error), PCA only *summarizes* the data.

Formative (PCA)



Reflective (EFA)



Exploratory Factor Analysis (EFA)

If $\mathbf{\Lambda}$ is not somehow constrained, latent variable variance is not identified. We can arbitrarily add rotation matrices \mathbf{T} and not change the decomposition:

$$\Sigma = \mathbf{\Lambda} \mathbf{T} \mathbf{T}^{-1} \mathbf{\Psi} \mathbf{T}^{-1\top} \mathbf{T}^{\top} \mathbf{\Lambda}^{\top} + \Theta$$

Can be seen as a different factor model with $\mathbf{\Lambda}^* = \mathbf{\Lambda} \mathbf{T}$ and $\mathbf{\Psi}^* = \mathbf{T}^{-1} \mathbf{\Psi} \mathbf{T}^{-1\top}$. To this end, in estimation one can assume uncorrelated factors, $\mathbf{\Psi} = \mathbf{I}$. Afterwards, rotation methods can be used to obtain simple structure for $\mathbf{\Lambda}$ while possibly allowing factors to correlate:

- orthogonal (varimax): axes remain orthogonal, independent
- oblique (promax/oblimin): axes become correlated

I always use promax.

Choosing the number of Factors is a bit more involved than PCA

- One method involves checking how many eigenvalues in $\mathbf{S} - \hat{\mathbf{\Theta}}$ are above 0
 - $\hat{\mathbf{\Theta}}$ is then estimated using a 1-factor model
- Parallel analysis takes sampling variation into account, and checks how many eigenvalues are statistically above what can be expected given an independence model

BFI example

```
library("psych")

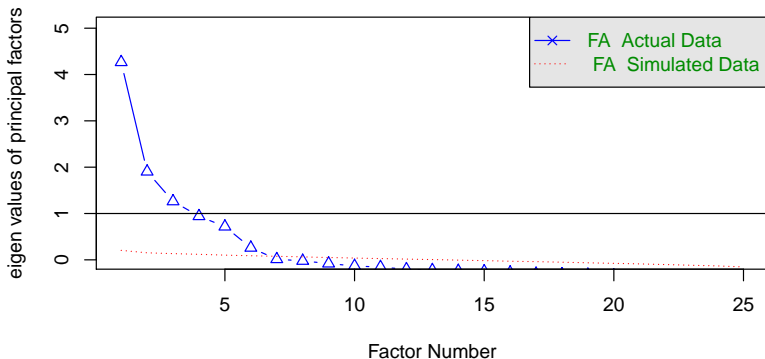
# Load data:
data(bfi)
bfiSub <- bfi[,1:25]

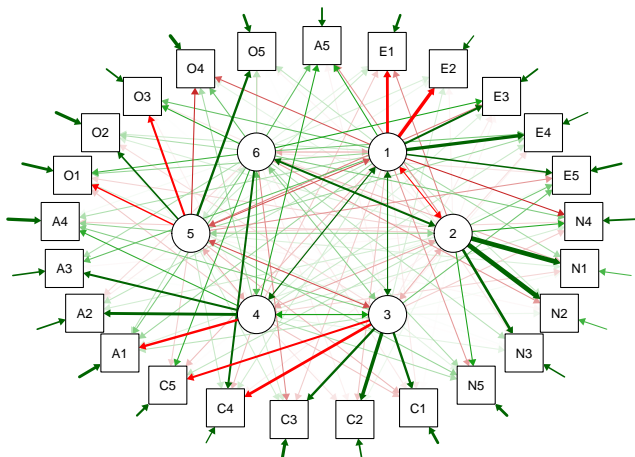
# Correlations:
corMat <- cor(bfiSub, use = "pairwise.complete.obs")
N <- nrow(bfiSub)
```

```
fa.parallel(corMat, N, fa = "fa")
```

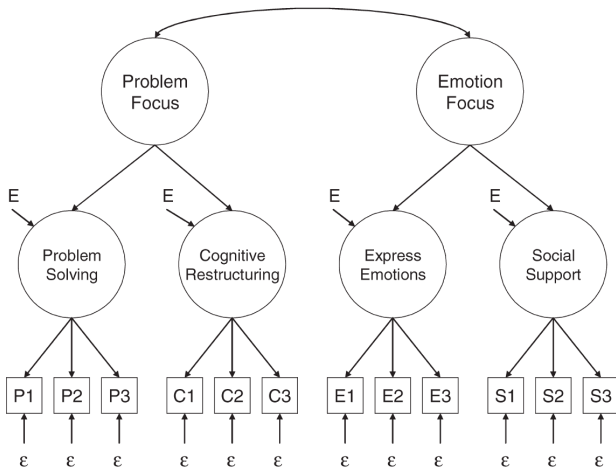
```
## Parallel analysis suggests that the number of factors = 6 and the
```

Parallel Analysis Scree Plots





Higher order models



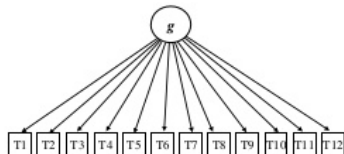
Higher order models

Mathematically, simply a second factor model on the latent variable variance–covariance matrix:

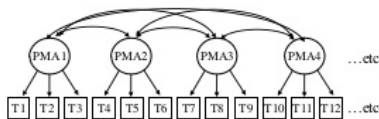
$$\Psi = \Lambda^* \Psi^* \Lambda^{*\top} + \Theta^*$$

Same rules of identification apply:

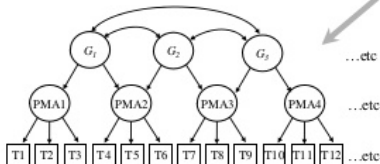
- The higher order factor must be scaled (one factor loading or the variance fixed to 1)
- The number of variances and covariances in Ψ must be at least as much as the number of parameters used to model Ψ



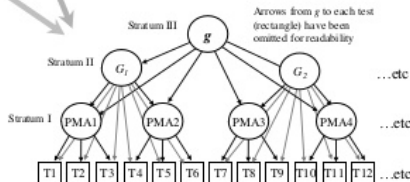
(1a) Spearman's general Factor model



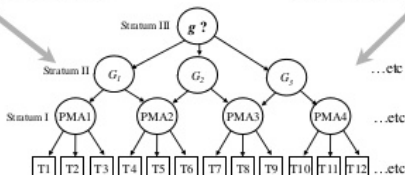
(1b) Thurston's Multiple Factor (Primary Mental Abilities) Model



(1c) Cattell-Horn *Gf-Gc* Hierarchical Model

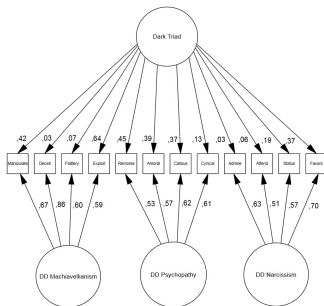


(1d) Carroll's Schmid-Leiman Hierarchical Three-Stratum Model



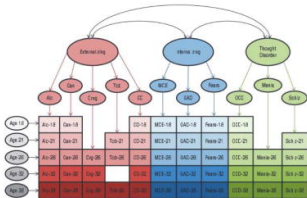
Note: Circles represent latent factors. Squares represent manifest measures (tests; T1...). Single-headed path arrows designate factor loadings. Double-headed arrows designate latent factor correlations

Bi-factor models

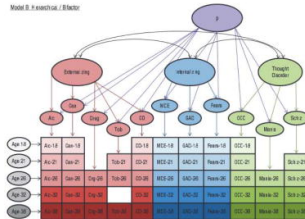


- Uncorrelated factors in combination with an uncorrelated bifactor
- Higher order model is nested in the bi-factor model
- Increasingly popular, but hard to interpret

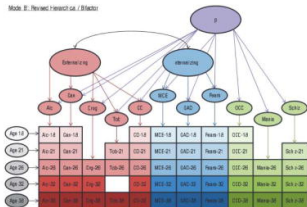
Model A: Correlated Factors



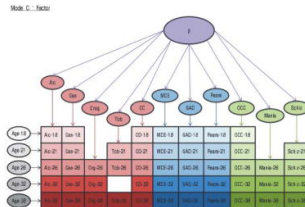
Model B: Hierarchical 2-Factor



Model B': Revised Hierarchical 2-Factor

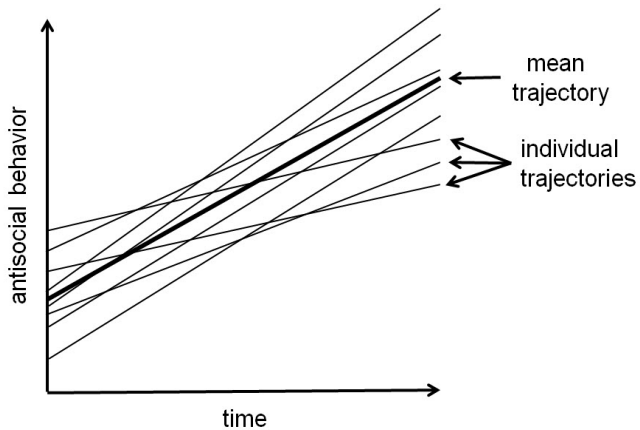


Model C: 1-Factor

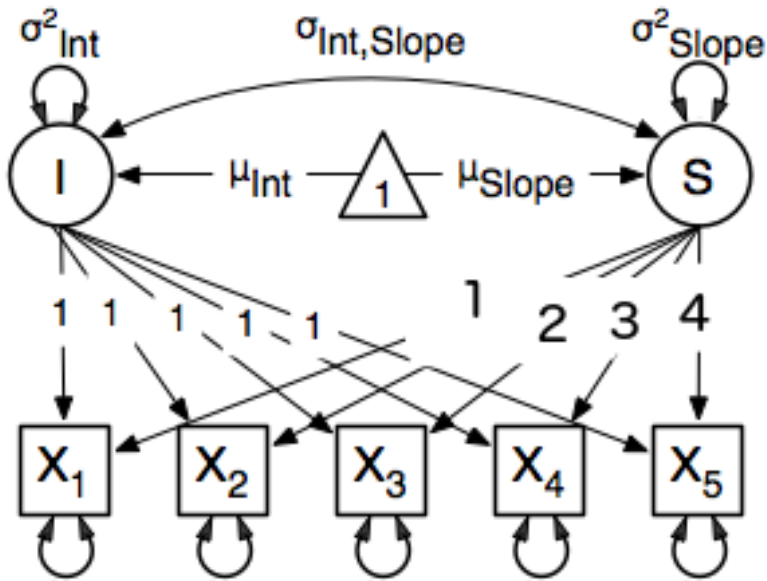


Caspi, A., Houts, R. M., Belsky, D. W., Goldman-Mellor, S. J., Harrington, H., Israel, S., ... & Moffitt, T. E. (2014). The p factor: one general psychopathology factor in the structure of psychiatric disorders?. *Clinical Psychological Science*, 2(2), 119-137.

Latent growth models



Latent growth models



- Exploratory factor analysis
- Missing data needs assumption of missing at random (MAR)
- Advanced CFA models:
 - Higher-order models
 - Bi-factor models
 - Latent growth models