

Structural Equation Modelling

Short course in Applied Psychometrics

*Peterhouse College,
Cambridge,
27-29 March 2012*

This course

The course is funded by the ESRC RDI and hosted by
The Psychometrics Centre



Tutors

Jon Heron, PhD (Bristol) jon.heron@bristol.ac.uk

Anna Brown, PhD (Cambridge) ab936@medschl.cam.ac.uk

Tim Croudace, PhD (Cambridge) tjc39@cam.ac.uk

	Day 1	Day 2	Day 3	
9:00-	Coffee on arrival	<u>Lec-6 – Special issues in CFA</u> Correlated errors Bi-factor modelling Method factors Multi-group CFA	<u>Lec-9 – SEM</u> Incorporating latent traits into path models.	9:00-
9:20-	Introductions + Aims of course			9:20-
9:40-	<u>Lec-1</u> Mplus modelling framework			9:40-
10:00-				10:00-
10:20-				10:20-
10:40-				10:40-
11:00-	Coffee	Coffee	Coffee	11:00-
11:20-	<u>Lec-2 – Regression models</u>	<u>Lec-7 – Path models 1</u> The basics / figures / Identification/ model fit/ equivalent models	<u>Examples 5 – SEM</u> EAS - SEM	11:20-
11:40-				11:40-
12:00-	<u>Examples 1</u> EAS - regression models		<u>Examples 3: SZ paper.</u>	Wrapping up, further reading and questions
12:20-		12:20-		
12:40-		12:40-		
13:00-	Lunch	Lunch	Lunch and depart	13:00-
13:20-				13:20-
13:40-				13:40-
14:00-	<u>Lec-3 - CFA with continuous variables</u>	<u>Lec-8 – Path models 2</u> Model refinement Direct and indirect effects Binary mediators - logit/probit		14:00-
14:20-				14:20-
14:40-	14:40-			
15:00-	<u>Lec-4 – EFA with continuous variables</u>			15:00-
15:20-				15:20-
15:40-				15:40-
16:00-	Coffee	Coffee		16:00-
16:20-	<u>Lec-5 - CFA and EFA with categorical variables</u>	<u>Examples 4</u> Path model using EAS		16:20-
16:40-				16:40-
17:00-	<u>Examples 2</u> EAS – CFA/EFA			17:00-
17:20-				17:20-
17:40-				17:40-

Learning objectives

- What are Measurement models
 - Constituent parts
 - Covariance structure
 - Identification
 - Difference between CFA and EFA
- How to deal with different types of observed variables
 - Continuous data
 - Binary/ordinal data
 - Polychoric correlations
- How to diagnose model misspecification
 - Residuals assessment
 - Model fit / Modification indices

Purposes of Factor analysis

- Assessment of dimensionality (i.e. how many latent variables underly your data...)
- Assessment of validity of items in questionnaires and surveys
- Scoring of respondents on latent variables
- Assessment of error of measurement
- Finding correlations among latent variables
- Answering specific scientific questions about relationship between observed and latent variables

Measurement models

- **Definition:** The mapping of measures onto theoretical constructs
- Measures – directly observed
- Theoretical constructs – unobserved (latent)



Latent variables

- It is something what we cannot observe/measure directly
- Also called **constructs** (especially in psychology), **factors**, **latent traits**
 - Intelligence, well-being, verbal ability, extraversion....
 - Continuous

Extraversion

- Just because a construct is named it does not mean it is correctly understood

Observed measures

- It is something what we can measure directly
- Also called **indicators, manifest variables, measures** and **proxies**
- Continuous
 - Height, head circumference, number of push-ups, time it takes to solve a maths problem
- Categorical
 - Responses on disagree – agree scales, never – always scales, school grades,.....

Height

Linear factor analysis model

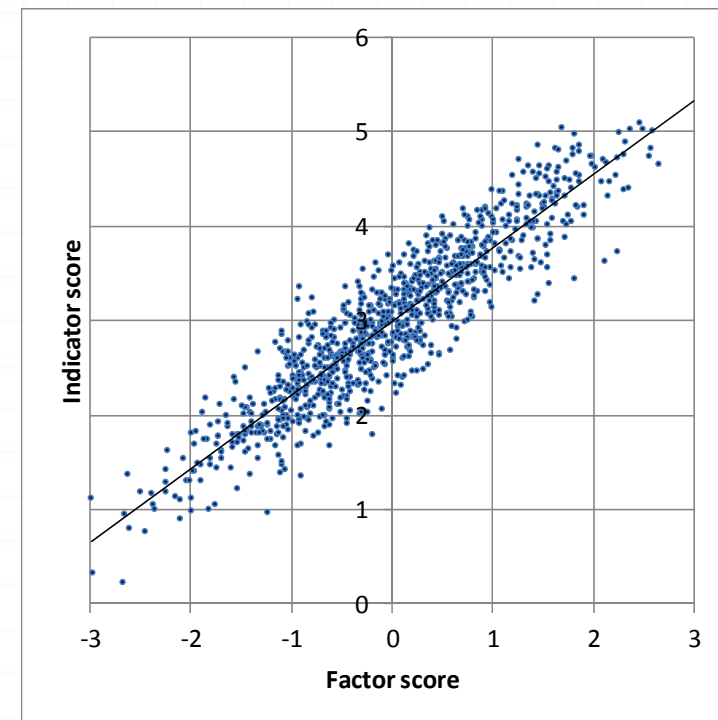
- We have observed random variables y_1, y_2, \dots, y_p with means m_1, m_2, \dots, m_p
- Each observed variable can be described as a function of a set of k common factors, and a unique factor

$$y_i = m_i + l_{i1} * F_1 + \dots + l_{ip} * F_k + e_i$$

- $k < p$
- We aim to describe the complexity in our p variables with a reduced number of variables (k)

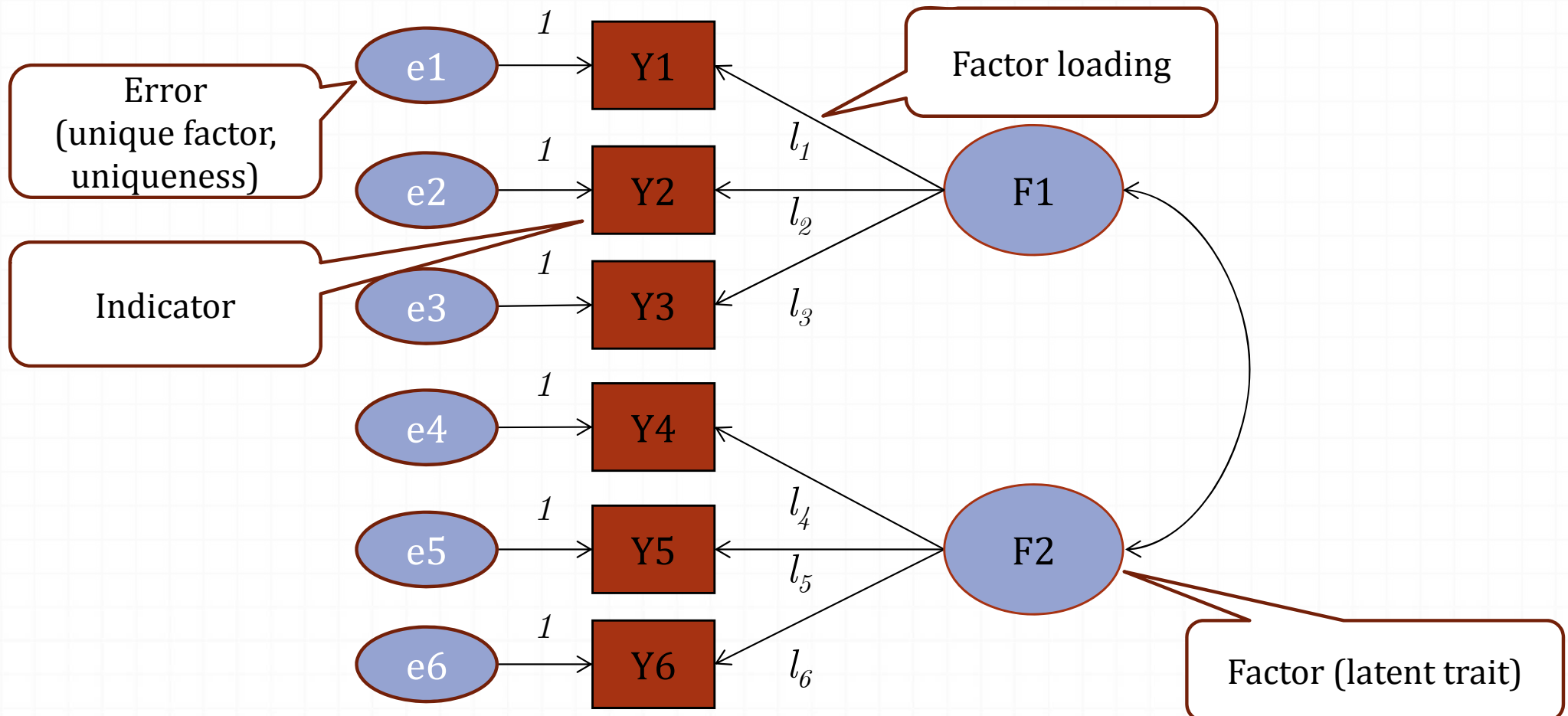
One factor model - example

- One common factor
- Observed measure can be described as a linear function of that factor
 - Mean becomes the regression line's *intercept*
 - Regression line's slope (*loading*) reflects the scale of change in the observed score when factor score changes 1 unit



CFA - Graphic representation

$$y_i = m_i + l_{i1} * F_1 + \dots + l_{ip} * F_k + e_i$$



Factor model – usual assumptions

- Factors and errors are independent

$$\text{cov}(F_d, e_i) = 0$$

- Errors are independent

$$\text{cov}(e_i, e_k) = 0$$

- Models with correlated errors will be considered later in the course

What needs to be estimated

- We know:
 - Variances of observed variables
 - Covariances of observed variables
- We don't know
 - Factor loadings l_{i1}, l_{i2}, \dots
 - Factor variances $var(F1), var(F2), \dots$
 - Factor covariances $cov(F1, F2), cov(F1, F3), \dots$
 - Error variances $var(e1), var(e2), \dots$
- If unique estimates exist for all parameters, the model is **identified**

Covariances

- Describe bivariate relationships

$$\text{cov}(X, Y) = \sum_{i=1}^N \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}$$

- May range from $-\infty$ to $+\infty$ (in theory)
- Covariance equals 0 for unrelated variables
- Difficult to say how „strong“ the relationship is without knowing the variables' variances

Variance-covariance matrix

- For our k observed variables: Y_1, Y_2, \dots, Y_k

$$\begin{bmatrix}
 \sum_{i=1}^N \frac{(y_{1i} - \bar{y}_1)(y_{1i} - \bar{y}_1)}{N} & \sum_{i=1}^N \frac{(y_{1i} - \bar{y}_1)(y_{2i} - \bar{y}_2)}{N} & \dots & \sum_{i=1}^N \frac{(y_{1i} - \bar{y}_1)(y_{ki} - \bar{y}_k)}{N} \\
 \sum_{i=1}^N \frac{(y_{2i} - \bar{y}_2)(y_{1i} - \bar{y}_1)}{N} & \ddots & & \vdots \\
 \vdots & & & \\
 \sum_{i=1}^N \frac{(y_{ki} - \bar{y}_k)(y_{1i} - \bar{y}_1)}{N} & \dots & & \sum_{i=1}^N \frac{(y_{ki} - \bar{y}_k)(y_{ki} - \bar{y}_k)}{N}
 \end{bmatrix}$$

Some properties of variances and covariances

○ For any constant c $\text{cov}(X, c) = 0$

○ Covariance is symmetrical $\text{cov}(X, Y) = \text{cov}(Y, X)$

○ Covariance of variable with itself is its variance

$$\text{cov}(X, X) = \text{var}(X)$$

○ Covariance of a sum

$$\text{cov}(X + Z, Y) = \text{cov}(X, Y) + \text{cov}(Z, Y)$$

○ Covariance between a variable and a product of a variable and a constant

$$\text{cov}(cX, Y) = c \text{cov}(X, Y)$$

Scale of latent variables

$$y_i = m_i + l_{i1} * F_1 + \dots + l_{ip} * F_k + e_i$$

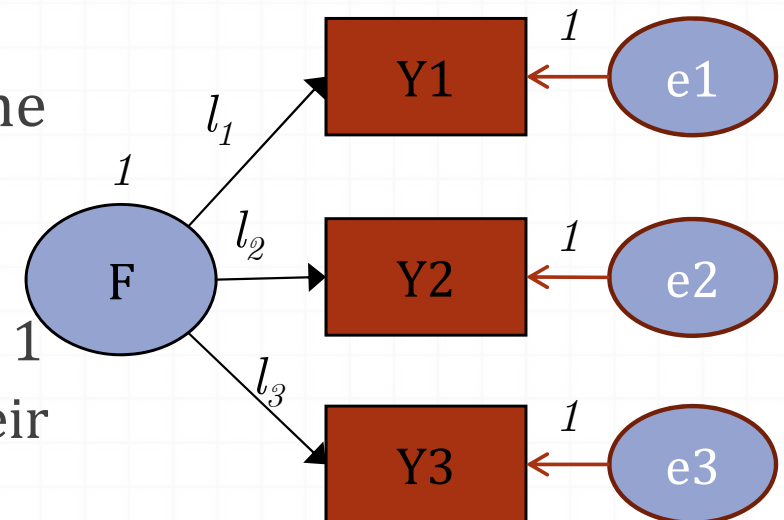
- o Latent variables have no scale (no origin, no unit)
- o We have to give them a scale (arbitrary)
 - o Set **mean = 0** for common and unique factors
 - o For **common factors** can either
 - o Set **variance = 1**
 - o Set scale according to the scale of **one of the measures**
 - o For **unique factors** we usually
 - o Set scale according to the scale of **the observed measure**

Setting scales in *Mplus*

- o Latent factors' scales need to be set either by
 1. Setting one factor loading – this is default in *Mplus*
`F1 BY y1 y2 y3; !read F1 BY y1@1 y2 y3;`
 - o Continuous observed variables have their own scale which they can pass to the latent factors
 2. Setting factor variance (to 1) and freeing the first factor loading:
`F1 BY y1* y2 y3; F1@1;`

Example - one factor model

- We have: 3 observed random variables
 - 3 variances and 3 covariances
- We want to explain this data through one common and 3 unique factors
- Model setup
 - Set the scale of F by fixing its variance to 1
 - Set the scale of uniquenesses by fixing their loadings to 1
- To estimate
 - 3 loadings, 3 unique variances
- This model is **just identified**



Covariance structure of the example 1-factor model

$$y_i = m_i + l_i * F + e_i$$

- o Scale for factor is set $\text{var}(F) = 1$
- o And the usual assumptions hold (independence of errors, and errors from the factor)

- o Then variance of any indicator (known)

$$\text{var}(y_i) = \text{var}(l_i * F) + 2 \text{cov}(l_i * F, e_i) + \text{var}(e_i) = l_i^2 + \text{var}(e_i)$$

- o Covariance of any 2 indicators (known)

$$\text{cov}(y_i, y_k) = \text{cov}(l_i * F + e_i, l_k * F + e_k) = \text{cov}(l_i * F, l_k * F) = l_i l_k$$

Equations for the example 1-factor model

- Then the 1-factor model with 3 indicators can be described by 6 equations:

$$\text{cov}(y_1, y_2) = l_1 l_2$$

$$\text{cov}(y_1, y_3) = l_1 l_3$$

$$\text{cov}(y_2, y_3) = l_2 l_3$$

$$\text{var}(y_1) = l_1^2 + \text{var}(e_1)$$

$$\text{var}(y_2) = l_2^2 + \text{var}(e_2)$$

$$\text{var}(y_3) = l_3^2 + \text{var}(e_3)$$

- There are 6 unknowns in the equations (left hand side in red are known values from a sample)
 - Loadings are estimated from the first 3 equations
 - Error variances are estimated from the last 3 equations

Example 2-factor model

$$y_i = m_i + l_{i1} * F_1 + 0 * F_2 + e_i \quad \text{for } i=1,2,3$$

$$y_k = m_k + 0 * F_1 + l_{k2} * F_2 + e_k \quad \text{for } k=4,5,6$$

Six indicators provide

6 variances +

15 covariances =

21 pieces of information

To be estimated

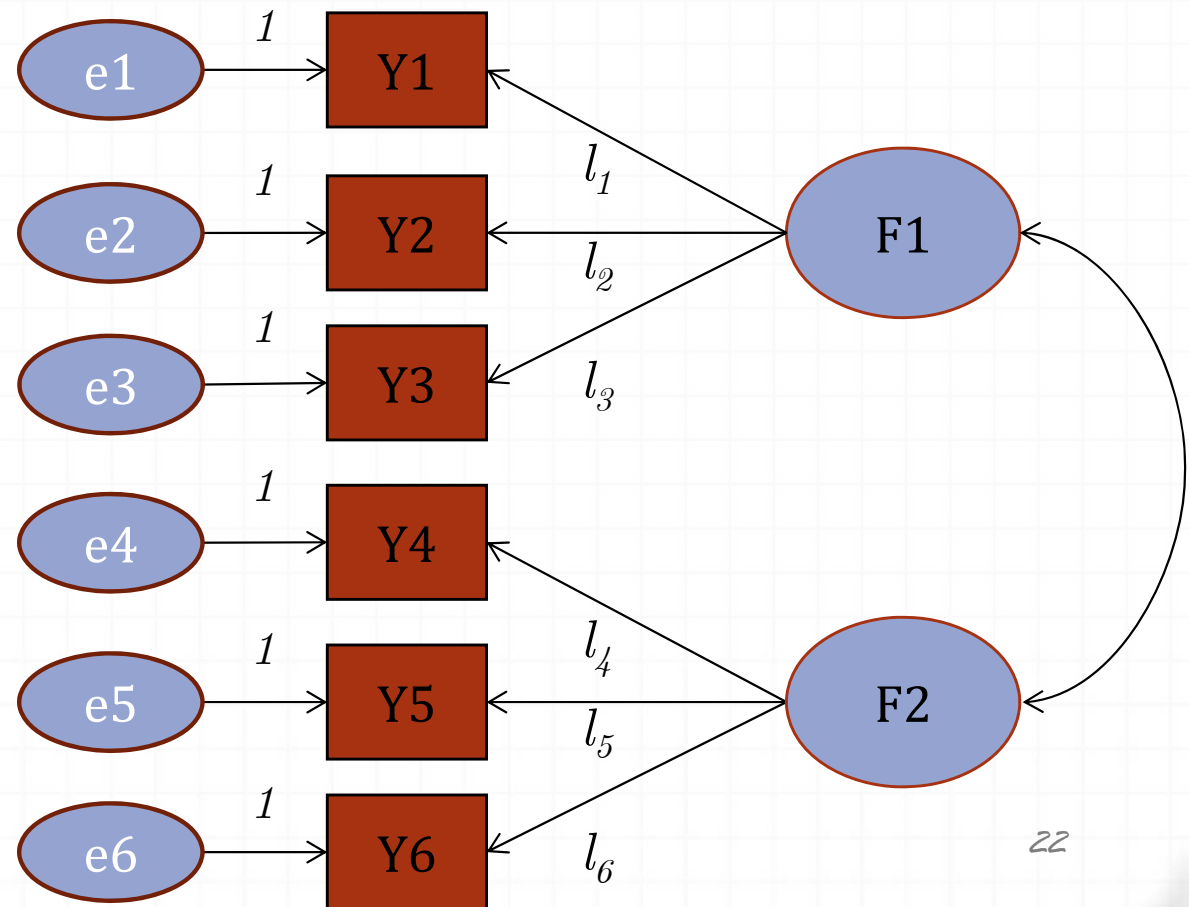
6 factor loadings +

6 error variances +

1 factor covariance =

13 pieces of information

This model is **over-identified**



Covariance structure of the example 2-factor model

- Scales for factors have been set $\text{var}(F_1) = \text{var}(F_2) = 1$
- And the usual assumptions hold (independence of errors, and errors from the factor)

- Then variance of any indicator (known)

$$\text{var}(y_i) = l_{i1}^2 + \text{var}(e_i)$$

- Covariance of any 2 indicators (known)

$$\begin{aligned}\text{cov}(y_i, y_k) &= \text{cov}(l_{i1} * F_1 + e_i, l_{k2} * F_2 + e_k) = \\ &= \text{cov}(l_{i1} * F_1, l_{k2} * F_2) = l_{i1} l_{k2} \text{cov}(F_1, F_2)\end{aligned}$$

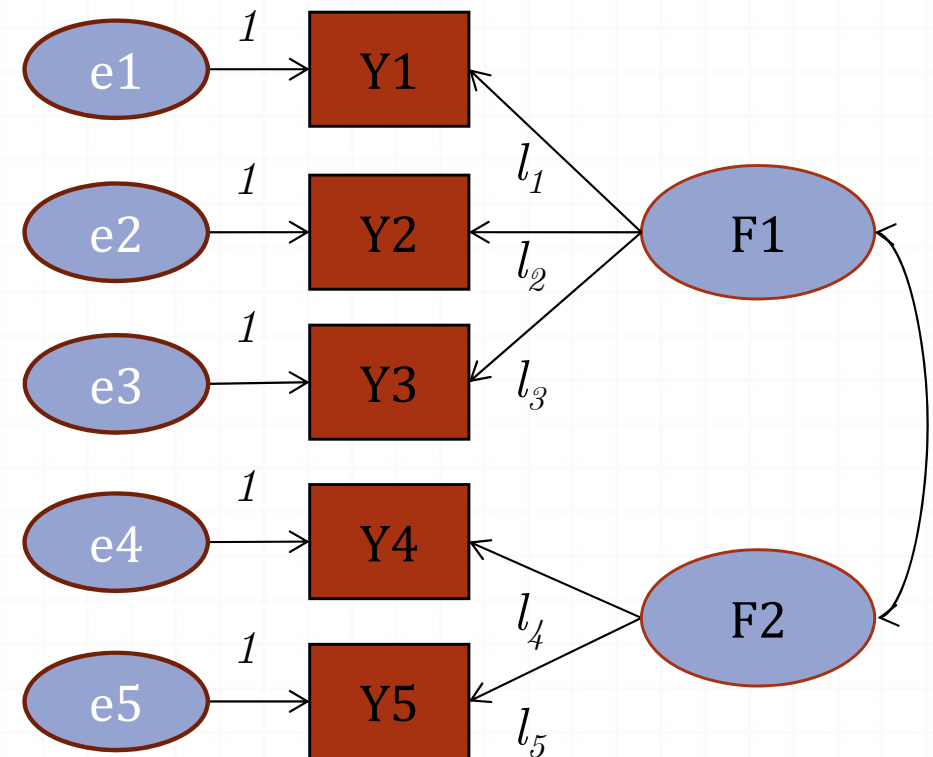
- This model is described by 21 equations, with only 13 unknowns

Identification rules for CFA

- **t-rule** (*necessary but not sufficient*) $t \leq \frac{1}{2} p(p+1)$
 - where t is the number of parameters to be estimated
- **General indicator rules**
 - One-factor model with 3 indicators is identified
 - Multifactor model is identified when it has (*sufficient but not necessary*)
 1. At least 2 indicators per factor
 2. Each indicator loads on only one factor
 3. Each row of factor covariance matrix has at least one non-zero off-diagonal element
 4. Errors are uncorrelated

Is this model identified?

- We will come across doublet-factors from time to time
 - Only 2 indicators
- What will happen if factors F1 and F2 are modelled as uncorrelated?
- Estimation of a doublet factor on its own requires breaking one covariance $\text{cov}(Y_4, Y_5)$ into 2 factor loadings
 - Not possible without constraints



Empirical identification

- o Model that is theoretically identified still might be empirically unidentified
- o Take an example of the 1-factor model with 3 indicators
 - o We know it is just identified setting $(\lambda_1) = 1$
 - o Now imagine that item 1 is not related to F in our dataset, i.e. $l_1 = 0$. Then two equations (see slide with [Equations](#) for this model) become uninformative

$$\text{cov}(y_1, y_2) = l_1 l_2 = 0$$

$$\text{cov}(y_1, y_3) = l_1 l_3 = 0$$

$$\text{cov}(y_2, y_3) = l_2 l_3$$

And it is not possible to resolve the equations uniquely in respect to l_2 and l_3

How CFA works in practice

1. You provide Mplus with covariance matrix (or raw data) for your sample
2. You specify hypothesis about underlying structure (how many factors and which items load on which factor)
3. Mplus will find parameters that conform to your hypothesis and maximise *likelihood** of observed data
4. Your real sample covariance matrix is compared to the covariance matrix based on estimated parameter values
5. If the difference is small enough, your data fits the model

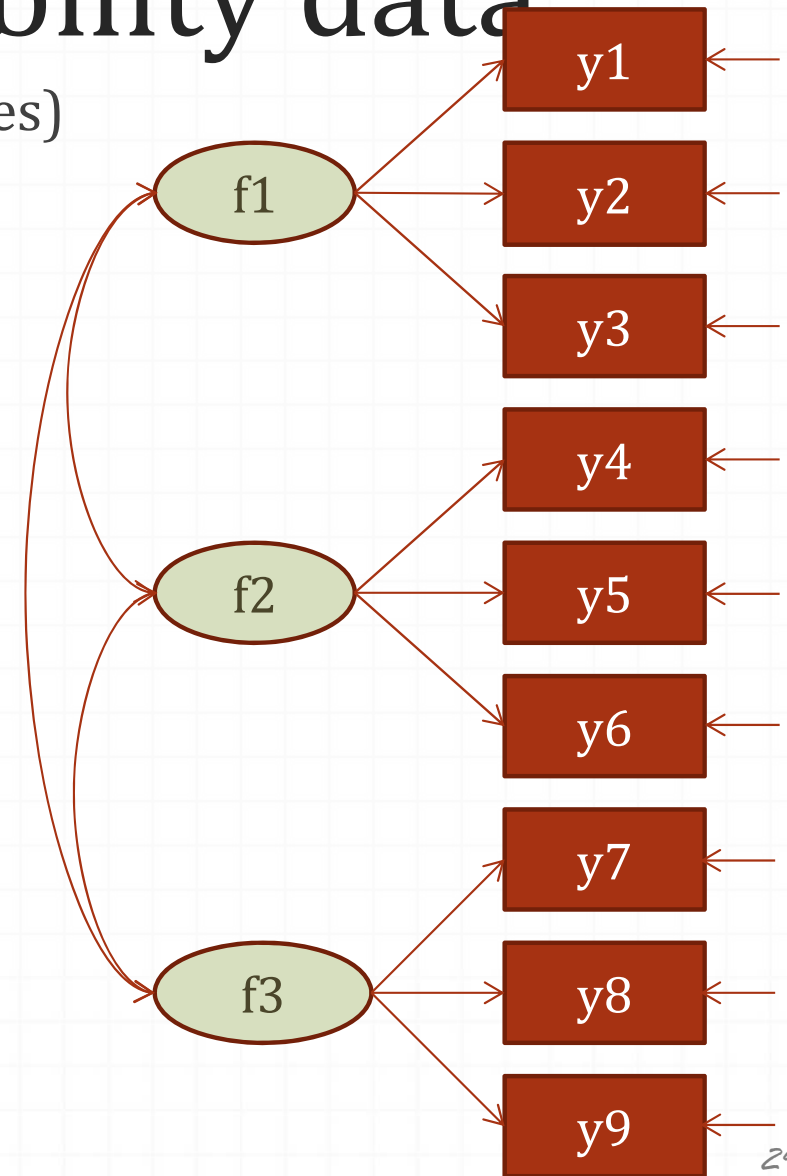
* assuming that ML estimator is used. Alternatively, other criteria are fulfilled.

Estimation

- Default estimator depends on type of analysis and measurement level of observed variables.
- For **continuous** variables default is ML
 - Minimizes **fit function** (related to discrepancy between sample covariances and those predicted by researcher model)
 - The model parameters obtained with this method maximize the likelihood of observing the available data if one were to collect data from the same population again
 - Assumes *multivariate normality*
- For non-normal continuous data, robust estimators are available
 - MLM is Mean corrected ML (Satorra/Bentler)
 - MLMV is Mean & Variance corrected ML (Muthen)
 - Both also produce (corrected) χ^2 -test and RMSEA

Example CFA: Thurstone's ability data

- o We have 9 subtests (continuous variables) measuring 3 Primary mental abilities
- o Verbal
 - o 1=sentences
 - o 2=vocabulary
 - o 3=sentence completion
- o Word fluency
 - o 4=first letters
 - o 5=four-letter words
 - o 6=suffixes
- o Reasoning
 - o 7=letter series
 - o 8=pedigrees
 - o 9=letter grouping



Thurstone's data – cont.

o We will analyse a correlation matrix (THUR.dat), $n=215$

1										1= sentences
.828	1									2= vocabulary
.776	.779	1								3= sentence completion
.439	.493	.460	1							4= first letters
.432	.464	.425	.674	1						5= four-letter words
.447	.489	.443	.590	.541	1					6= suffixes
.447	.432	.401	.381	.402	.288	1				7= letter series
.541	.537	.534	.350	.367	.320	.555	1			8= pedigrees
.380	.358	.359	.424	.446	.325	.598	.452	1		9= letter grouping

What needs to be estimated

○ Factor loadings

$$\begin{pmatrix} l_1 & 0 & 0 \\ l_2 & 0 & 0 \\ l_3 & 0 & 0 \\ 0 & l_4 & 0 \\ 0 & l_5 & 0 \\ 0 & l_6 & 0 \\ 0 & 0 & l_7 \\ 0 & 0 & l_8 \\ 0 & 0 & l_9 \end{pmatrix}$$

○ Factor covariances

$$\begin{pmatrix} \text{var}(F_1) & & \\ \text{cov}(F_1, F_2) & \text{var}(F_2) & \\ \text{cov}(F_1, F_3) & \text{cov}(F_2, F_3) & \text{var}(F_3) \end{pmatrix}$$

○ Error variances

$$\begin{pmatrix} \text{var}(e_1) & & & & \\ 0 & \text{var}(e_2) & & & \\ \vdots & & \ddots & & \\ 0 & 0 & \cdots & \text{var}(e_9) \end{pmatrix}$$

Thurstone's data – CFA syntax

```
TITLE: CFA of Thurstone correlation matrix
DATA: FILE IS THUR.dat;
      TYPE IS CORRELATION;
      NOBSERVATIONS = 215;
VARIABLE: NAMES ARE subtest1-subtest9;
ANALYSIS: !defaults are ok; maximum likelihood
MODEL:
test1 BY subtest1-subtest3*;
test2 BY subtest4-subtest6*;
test3 BY subtest7-subtest9*;
test1-test3@1;
OUTPUT: RES;
PLOT: TYPE=PLOT2;
```


Thurstone's data - results

		Estimate	S.E.	Est./S.E.	Two-Tailed P-Value
TEST1	BY				
	SUBTEST1	0.903	0.054	16.805	0.000
	SUBTEST2	0.912	0.053	17.084	0.000
	SUBTEST3	0.854	0.056	15.388	0.000
TEST2	BY				
	SUBTEST4	0.834	0.060	13.847	0.000
	SUBTEST5	0.795	0.061	12.998	0.000
	SUBTEST6	0.701	0.064	11.012	0.000
TEST3	BY				
	SUBTEST7	0.779	0.064	12.231	0.000
	SUBTEST8	0.718	0.065	11.050	0.000
	SUBTEST9	0.702	0.065	10.729	0.000
TEST2	WITH				
	TEST1	0.643	0.050	12.815	0.000
TEST3	WITH				
	TEST1	0.670	0.051	13.215	0.000
	TEST2	0.637	0.058	10.951	0.000

SE of estimates
are of order
0.07 or below
 $1/\sqrt{215}=0.068$

Thurstone's data – residuals

o Model explains most correlations well

	1	2	3	4	5	6	7	8	9
SUBTEST1	0.000								
SUBTEST2	0.001	0.000							
SUBTEST3	0.001	-0.003	0.000						
SUBTEST4	-0.047	0.002	0.000	0.000					
SUBTEST5	-0.031	-0.004	-0.014	0.008	0.000				
SUBTEST6	0.038	0.076	0.056	0.003	-0.019	0.000			
SUBTEST7	-0.026	-0.046	-0.047	-0.035	0.005	-0.061	0.000		
SUBTEST8	0.104	0.096	0.120	-0.033	0.001	-0.002	-0.007	0.000	
SUBTEST9	-0.046	-0.072	-0.044	0.049	0.088	0.010	0.048	-0.054	0.000

Thurstone's data - model fit

Chi-Square Test of Model Fit

Value	38.737
Degrees of Freedom	24
P-Value	0.0291

CFI	0.986
-----	-------

RMSEA (Root Mean Square Error Of Approximation)

Estimate	0.053
90 Percent C.I.	0.017 0.083

SRMR (Standardized Root Mean Square Residual)

Value	0.044
-------	-------

o What does it all mean???

Minimum Fit Function Chi-Square = 283.78 (P = 0.0)
 Normal Theory Weighted Least Squares Chi-Square = 243.60 (P = 0.0)
 Satorra-Bentler Scaled Chi-Square = 29.05 (P = 0.010)
 Chi-Square Corrected for Non-Normality = 35.79 (P = 0.0011)
 Estimated Non-centrality Parameter (NCP) = 15.05
 90 Percent Confidence Interval for NCP = (3.33 ; 34.52)
 Minimum Fit Function Value = 0.60
 Population Discrepancy Function Value (F0) = 0.032
 90 Percent Confidence Interval for F0 = (0.0071 ; 0.073)
 Root Mean Square Error of Approximation (RMSEA) = 0.048
 90 Percent Confidence Interval for RMSEA = (0.022 ; 0.072)
 P-Value for Test of Close Fit (RMSEA < 0.05) = 0.52
 Expected Cross-Validation Index (ECVI) = 0.12
 90 Percent Confidence Interval for ECVI = (0.096 ; 0.16)
 ECVI for Saturated Model = 0.12
 ECVI for Independence Model = 11.74
 Chi-Square for Independence Model with 21 Degrees of Freedom = 5514.61
 Independence AIC = 5528.61
 Model AIC = 57.05
 Saturated AIC = 56.00
 Independence CAIC = 5564.71
 Model CAIC = 129.25
 Saturated CAIC = 200.40
 Normed Fit Index (NFI) = 0.99
 Non-Normed Fit Index (NNFI) = 1.00
 Parsimony Normed Fit Index (PNFI) = 0.66
 Comparative Fit Index (CFI) = 1.00
 Incremental Fit Index (IFI) = 1.00
 Relative Fit Index (RFI) = 0.99
 Critical N (CN) = 473.52
 Root Mean Square Residual (RMR) = 0.038
 Standardized RMR = 0.038
 Goodness of Fit Index (GFI) = 0.87
 Adjusted Goodness of Fit Index (AGFI) = 0.74
 Parsimony Goodness of Fit Index (PGFI) = 0.44

Goodness of fit

- Goodness of fit statistics are based on different ideas (e.g. summarizing elements in residual matrix, information theory, etc.)
- Some of them are known to favour certain types of model
- Fortunately Mplus provides only few of them (the ones that are known to provide good information about model fit)



Goodness of fit

Chi-square and log-likelihood

- **Null hypothesis:** the population covariance matrix **is equal** to the model-based estimated covariance matrix
 - “Accept-reject” hypothesis
 - Setting significance level very low works FOR researcher’s model
- **Chi-square** is widely used, although it has some undesirable properties
 - χ^2 as well $\Delta \chi^2$ are affected by sample size and model complexity (larger samples and more complex models tend to be rejected)
- **Log-likelihood** value is used to compare nested models
 - 2 x loglikelihood follows chi-square distribution with df equal to difference in number of estimated parameters

Goodness of fit

Information criteria

o Akaike Information Criterion (AIC)

$$\text{AIC} = \chi^2 + p(p - 1) - 2df$$

- o Note that because $p(p - 1)/2 - k = df$, $p(p - 1) - 2df$ is double the number of free parameters in the model
 - o So every free parameter pays a penalty of 2
 - o It is meaningful only when two different models are estimated.
 - o Lower values indicate a better fit and so the model with the lowest AIC is the best fitting model.
- ### o Sample-size adjusted Bayesian Information Criterion (BIC)

$$\chi^2 + \ln(N)[k(k - 1)/2 - df]$$

Goodness of fit

CFI

- Compare your model with baseline model
 - all observed variables are uncorrelated (terrible model!)
- **Comparative Fit Index (CFI)**

$$d = \chi^2 - df$$

$$\text{CFA} = \frac{d(\text{Null Model}) - d(\text{Proposed Model})}{d(\text{Null Model})}$$

- Ranges from 0 to 1, the higher the better, often recommended cutoff 0.95

Goodness of fit

Error of Approximation

- Root-mean-square Error of Approximation (**RMSEA**)

$$RMSEA = \sqrt{\frac{\chi^2 - df}{df (N - 1)}}$$

- Only RMSEA is not affected by model complexity (Cheung and Rensvold, 2002)
- RMSEA has a known distribution (non-central chi-square) and therefore confidence intervals can be computed
- Rules of thumb
 - MacCallum, Browne and Sugawara (1996) have used 0.01, 0.05, and 0.08 to indicate excellent, good, and mediocre fit respectively; RMSEA > 1 is considered poor fit

Goodness of fit

Residual Based Fit Indices

- o Measure average differences between sample and estimated population covariance (correlation) matrix
- o
- o Standardised root mean square residual (**SRMR**)
 - o ranges from 0 to 1
 - o the smaller the better, recommended cutoff 0.08

Goodness of fit - recommendations

- The goodness of fit indices address **global** fit
- Some argue that instead relying on these indices, the researcher should **always** locate the source of specification error
 - Check residual matrix for areas of local misfit
- Do not make your decisions on the basis of one fit index
 - There is always at least one fit index that shows good fit of your model (McDonald, 1999)
- When the sample size is big or when the model is complex, use other statistics than chi-square

Modification Indices

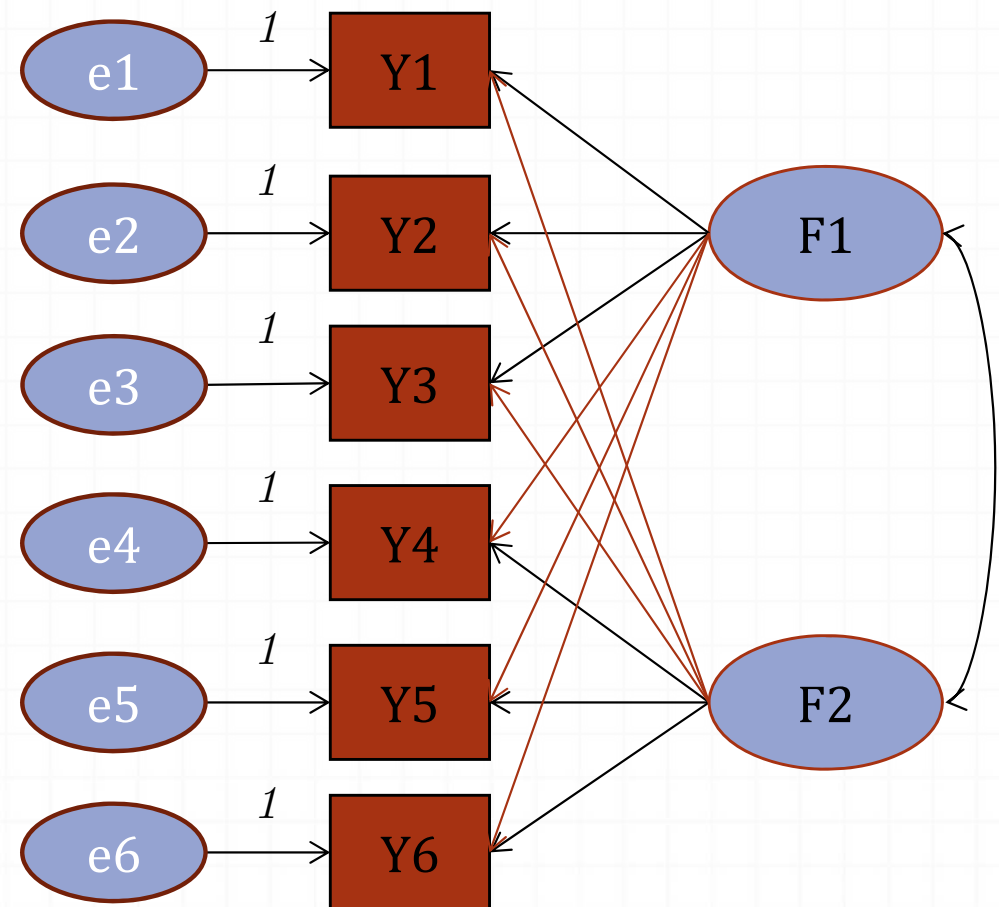
- Modification index (M.I.) is the value by which **chi-square** will drop if the parameter currently fixed or constrained was freely estimated
- Useful to guide modification of the model
- To request modification indices
OUTPUT: MOD (<*min.value*>);
- Only univariate MIs are available in *Mplus*
 - E.P.C. is expected parameter change index
 - Expected value of the parameter if it was freely estimated

Exploratory factor analysis

Learning to perform and interpret EFA

EFA model

- Purpose
 - Exploring a structure
 - Relaxing very strict assumptions imposed by CFA
- EFA can be considered as a special case of CFA model (or the other way around)
- Factors are indicated by **all** observed variables
- Factors can co-vary or not



EFA Identification 1

- In the single-factor case, the model is *identified*
- In the more general case of 2 or more factors, the system of equations describing the variables through common factors does not have a unique solution
 - There are infinite number of models that fit the data equally well
 - Further constraints are required
- Exchange of factor loadings while unique variances are identified and unchanging is called *rotation problem*
 - Resolved by assigning arbitrary loadings and then “rotating” them to approximate a given model (on this later)

Covariance structure of an EFA model

- Take a two-factor model as an example

$$y_i = m_i + l_{i1} * F_1 + l_{i2} * F_2 + e_i$$

- Scales for factors are set $\text{var}(F_1) = \text{var}(F_2) = 1$

- Then variance of any indicator

$$\text{var}(y_i) = l_{i1}^2 + l_{i2}^2 + 2l_{i1}l_{i2} \text{cov}(F_1, F_2) + \text{var}(e_i)$$

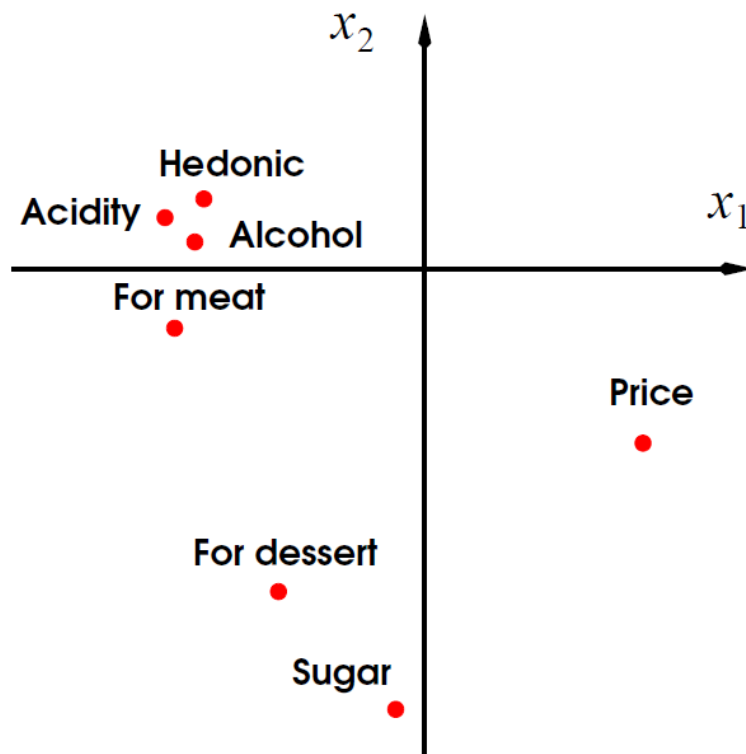
- Covariance of any 2 indicators

$$\text{cov}(y_i, y_k) = l_{i1}l_{k1} + l_{i2}l_{k2} + (l_{i1}l_{k2} + l_{k1}l_{i2})\text{cov}(F_1, F_2)$$

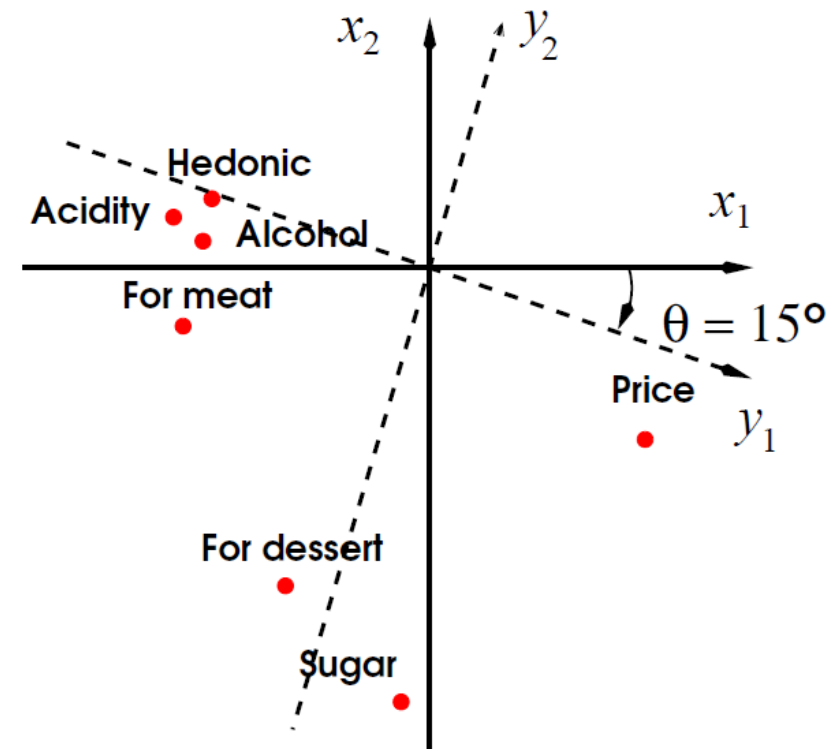
- The loadings are not identified

Rotation- example

Original factor loadings



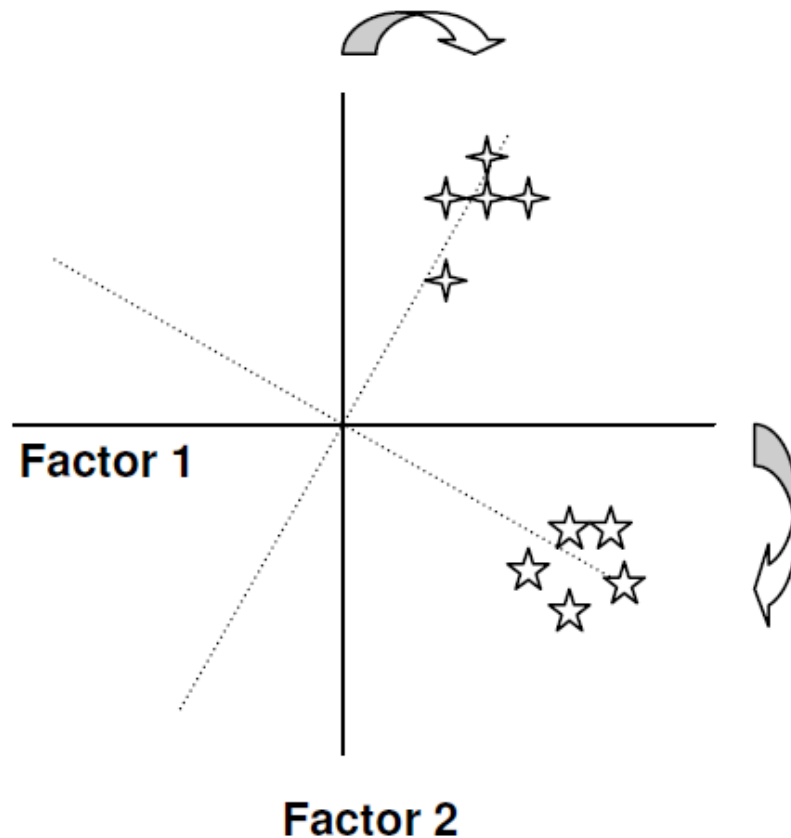
Orthogonal rotation



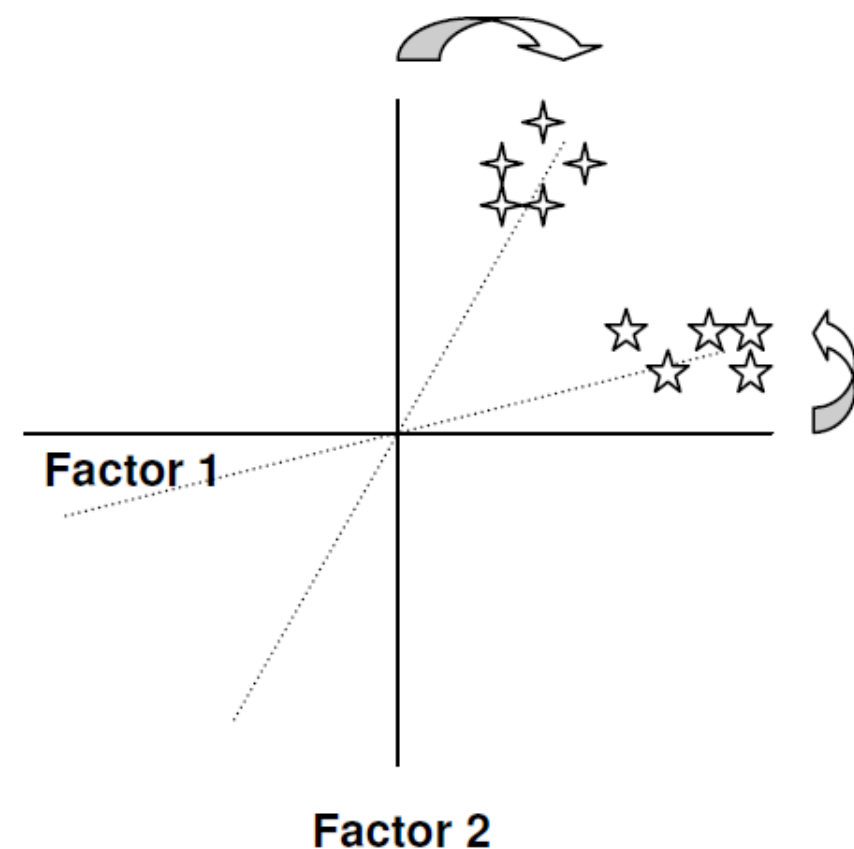
Abdi, H. (2003). Factor Rotations in Factor Analyses. In Lewis-Beck M., Bryman, A., Futing T. (Eds.), Encyclopedia of Social Sciences Research Methods. Thousand Oaks (CA): Sage.

Orthogonal versus oblique rotation

○ Orthogonal rotation



○ Oblique rotation



Independent clusters

- Item or test that indicates only 1 factor is called *factorially simple*
- Item or test that indicates 2 or more factor is called *factorially complex*
- *Independent clusters factor model* – every variable is an indicator for only 1 factor (every variable is factorially simple)

Rotation 1

- Rotation is a transformation of parameters to approximate an independent cluster solution (usually)
- Factors are uncorrelated (*orthogonal rotation*) or correlated (*oblique rotation*)
- McDonald (Test Theory, 1999) shows convincingly why oblique rotations are to be preferred
 - They avoid identification problems which will create “doublets” factors
 - For most applications correlated factors are more conceptually sound
 - Even if factors are found to be uncorrelated in one population, they might be correlated in another

Rotation 2

- o Many rotation algorithms are available in *Mplus*
- o For orthogonal rotations
 - o There are just rotated loadings to interpret
- o For oblique rotations
 - o There is a **pattern matrix** (like coefficients in multiple regression - correlations between indicators and the factor with other indicators partialled out)
 - o There is also a **structure matrix** (correlations between indicators and the factor)
 - o Correlations between the factors

EFA Identification 2

- o Another form of lack of identifiability
- o Joint indeterminacy of factor loadings and unique variances – hidden *doublet factors*
 - o Happens because for just two tests, $\sigma_{12} = \lambda_1 \lambda_2$ cannot be solved uniquely for λ_1 and λ_2
 - o In EFA with uncorrelated factors this cannot be resolved and is hidden by the analysis

Conducting EFA in practice

- o Model identification considerations
 - o Choice of rotation
 - o Checking the standard errors (ensuring identification)
 - o Checking the fit and the residuals
-
- o Main reference: McDonald, R. (1999). *Test Theory*. Lawrence Erlbaum.

EFA command in *Mplus*

ANALYSIS:

TYPE = EFA # #;

ROTATION = **GEOMIN**; ! (OBLIQUE) - default or (ORTHOGONAL)

QUARTIMIN !oblique only

CF-VARIMAX

CF-QUARTIMAX

CF-EQUAMAX

CF-PARSIMAX

CF-FACPARSIM

CRAWFER

OBLIMIN

PROMAX !oblique only

VARIMAX !orthogonal only

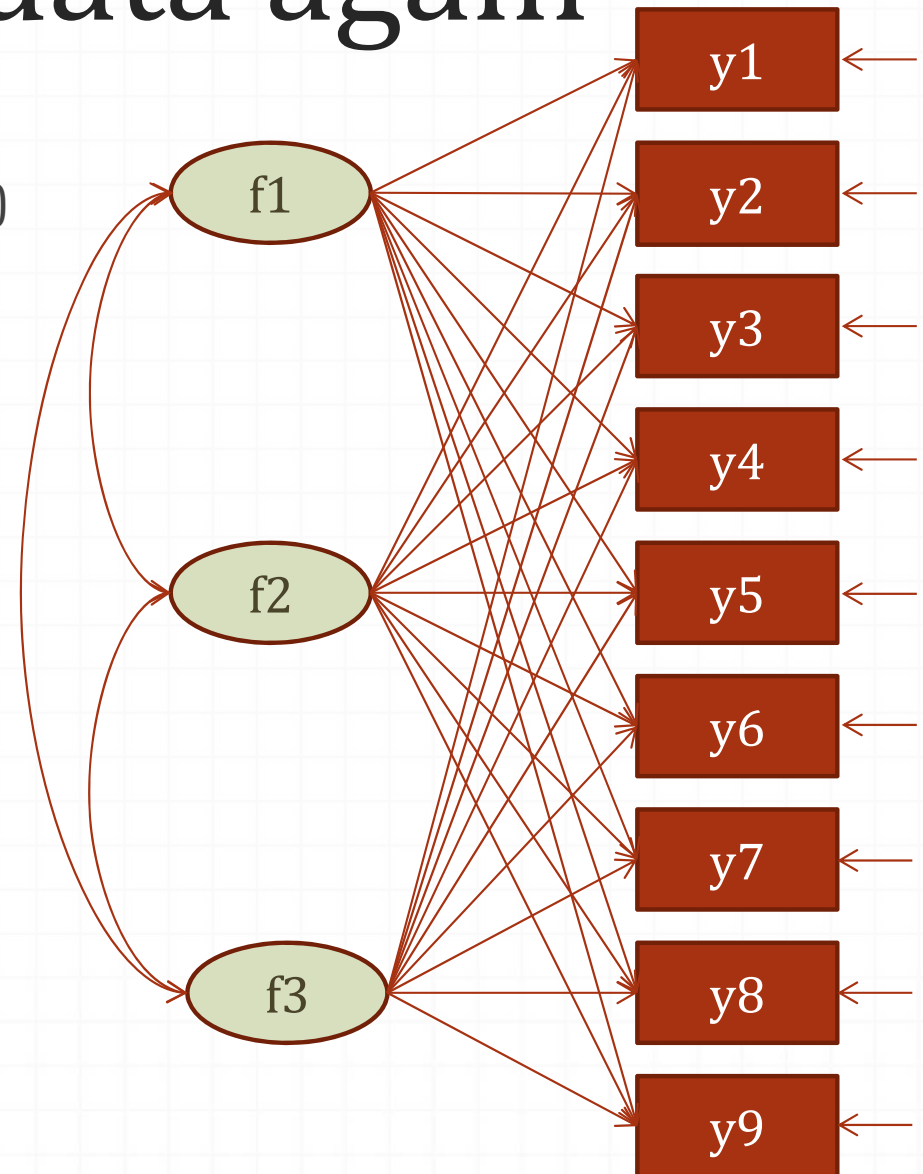
TARGET

Checking the standard errors

- For an identified model, SE should be approximately equal $1 / \sqrt{n}$
- If so, it is safe to proceed with the exploratory analysis
- If not, it might indicate an indeterminacy with doublet factors

Thurstone's data again

- o We have 9 subtests (continuous variables) measuring 3 Primary mental abilities
- o **Verbal**
 - o 1=sentences
 - o 2=vocabulary
 - o 3=sentence completion
- o **Word fluency**
 - o 4=first letters
 - o 5=four-letter words
 - o 6=suffixes
- o **Reasoning**
 - o 7=letter series
 - o 8=pedigrees
 - o 9=letter grouping



Thurstone's data – syntax for EFA

TITLE: EFA of Thurstone correlation matrix of Primary mental abilities

DATA: FILE IS THUR.dat;

TYPE IS CORRELATION;

NOBSERVATIONS = 215;

VARIABLE: NAMES ARE subtest1-subtest9;

ANALYSIS:

TYPE IS EFA 1 3; !we will fit 1, 2 and 3 factor models

ROTATION=CF-VARIMAX (ORTHOGONAL);

!ROTATION=CF-VARIMAX (OBLIQUE);

OUTPUT: RESIDUALS; !optional, very useful in model assessment

PLOT: TYPE = PLOT2; !optional, will produce a scree plot

Thurstone's data – Eigenvalues

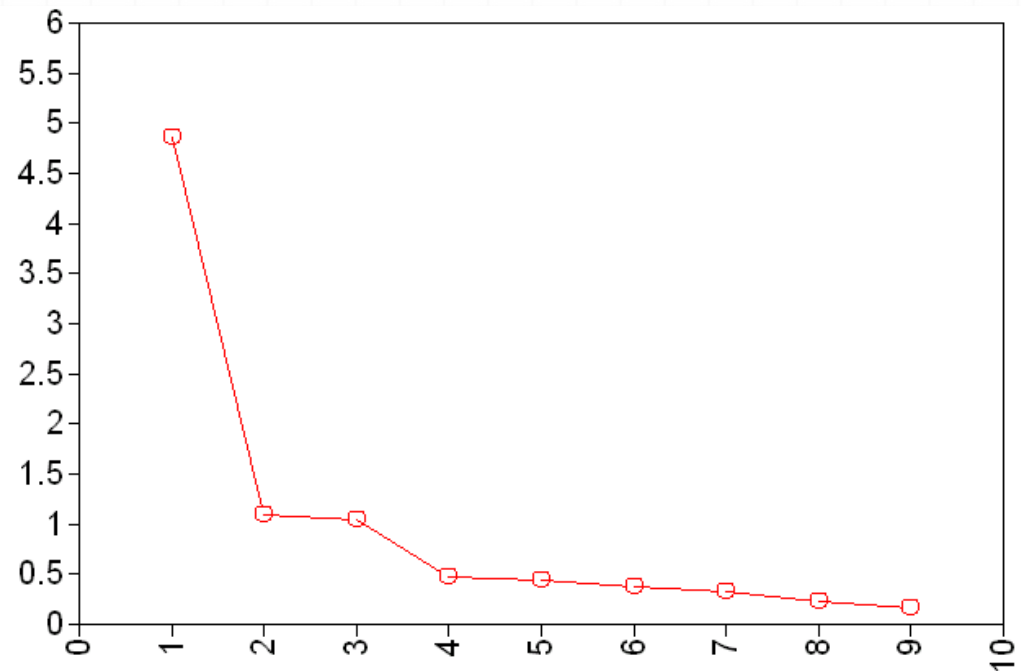
EIGENVALUES FOR SAMPLE CORRELATION MATRIX

1	2	3	4	5	6	7	8	9
4.851	1.090	1.038	0.475	0.448	0.375	0.321	0.234	0.168

o Scree plot

o **PLOT** command;

o TYPE=PLOT3;



Thurstone's data – residuals

- In the 2-factor model correlations between the last 3 subtests are not explained well

	SUBTEST6	SUBTEST7	SUBTEST8	SUBTEST9
SUBTEST6	0.000			
SUBTEST7	-0.086	0.000		
SUBTEST8	-0.048	0.217	0.000	
SUBTEST9	-0.062	0.284	0.143	0.000

- 3-factor model has near-0 residuals
- We will proceed with 3 factors for this data

Thurstone's data – model fit

	1 factor	2 factors	3 factors
Chi square	236.848	86.112	2.944
df	27	19	12
CFI	.806	.938	1
RMSEA	.190	.128	0

- 3 factor model is over fitting but 2 factor model is clearly not acceptable
- Check standard errors – are they of magnitude $1 / \sqrt{n}$ (is the model identified?)
 - Sample size is $n=215$, so SE should be of order 0.07

Thurstone's data – orthogonal rotated loadings

	1	2	3
SUBTEST1	0.858	0.196	0.223
SUBTEST2	0.854	0.270	0.180
SUBTEST3	0.800	0.240	0.187
SUBTEST4	0.287	0.782	0.197
SUBTEST5	0.269	0.698	0.261
SUBTEST6	0.358	0.598	0.103
SUBTEST7	0.277	0.185	0.779
SUBTEST8	0.478	0.151	0.503
SUBTEST9	0.200	0.317	0.622

- o Factor loadings are largely in line with expectations, however, there are many non-zero loadings

Thurstone's data – oblique rotated loadings

	1	2	3
SUBTEST1	0.824	0.044	0.121
SUBTEST2	0.811	0.139	0.058
SUBTEST3	0.758	0.111	0.078
SUBTEST4	0.025	0.817	0.053
SUBTEST5	0.011	0.709	0.145
SUBTEST6	0.187	0.614	-0.031
SUBTEST7	0.016	-0.003	0.842
SUBTEST8	0.332	-0.012	0.501
SUBTEST9	-0.061	0.198	0.643

o Factor loadings are much closer to an independent clusters solution

Thurstone's data – Factor correlations

o Factor correlations in the oblique solution

	1	2	3
1	1.000		
2	0.463	1.000	
3	0.455	0.464	1.000

o We would expect mental abilities to be correlated

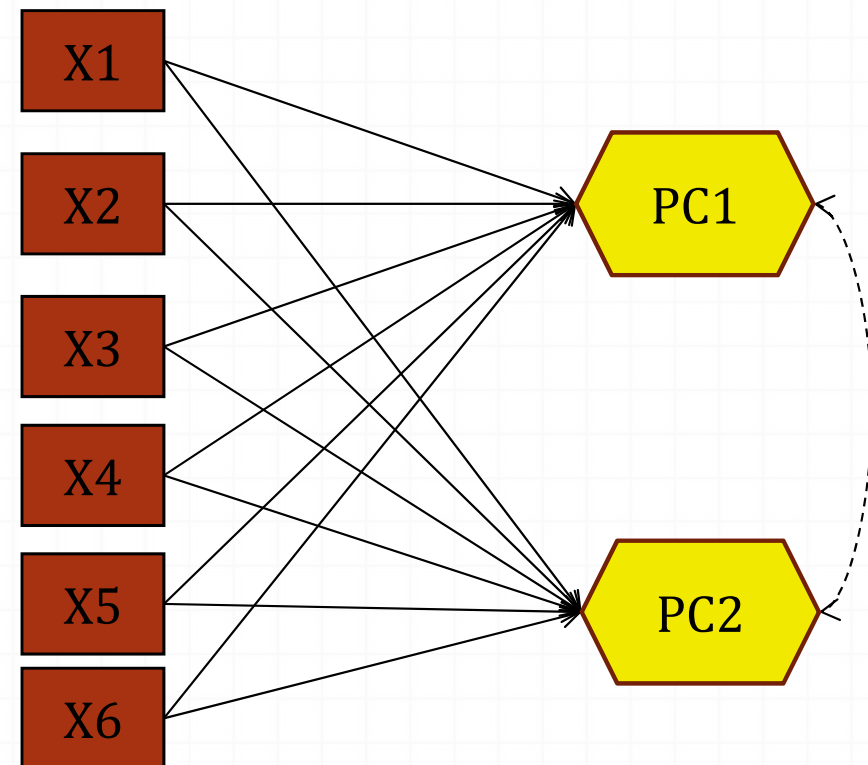
o Happy with the solution with 3 correlated factors

General notes on factor analysis

- Confirmatory factor analysis (**CFA**) – theory driven, more parsimonious and scientifically more sound methodology for finding underlying factors
- Exploratory factor analysis (**EFA**) – data driven automated searching engine for finding underlying factors
- Principal component analysis (**PCA**) – many think of it as one type of factor analysis, but PCA is conceptually different!

Principal component analysis (PCA)

- PCs are conveniently weighted sum-scores
- Constructs are casually determined by the observations
 - **Formative** measurement
 - EFA and CFA are **reflective** measurement
- Unique variances are missing (thus we do not account for measurement error)



PCA versus EFA

- PCA and EFA may look similar and in practice may look like giving similar results. But the principal components (from PCA analysis) and factors (from EFA analysis) have very different interpretations
- Use **EFA** when you are interested in *making statements about the factors that are responsible for a set of observed responses*
- Use **PCA** when you are simply interested in *performing data reduction*.

CFA and EFA with categorical variables

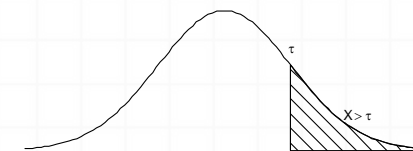
Categorical measures

- o Many observed variables are categorical, including test items
 - o Ability tests most often have *binary* responses (correct – incorrect)
 - o Questionnaires that employ rating scales most often have ordered categorical (*ordinal*) responses (often 3, 4 or 5)
 - o Rating scales can be symmetrical (agree-disagree) and not (never-always)
 - o Many rating categories (for instance, 9) sometimes allow treating data as continuous
- o Already learnt from *regression* that relationships between categorical variables and continuous factors are **non-linear**

Binary indicators and response tendencies

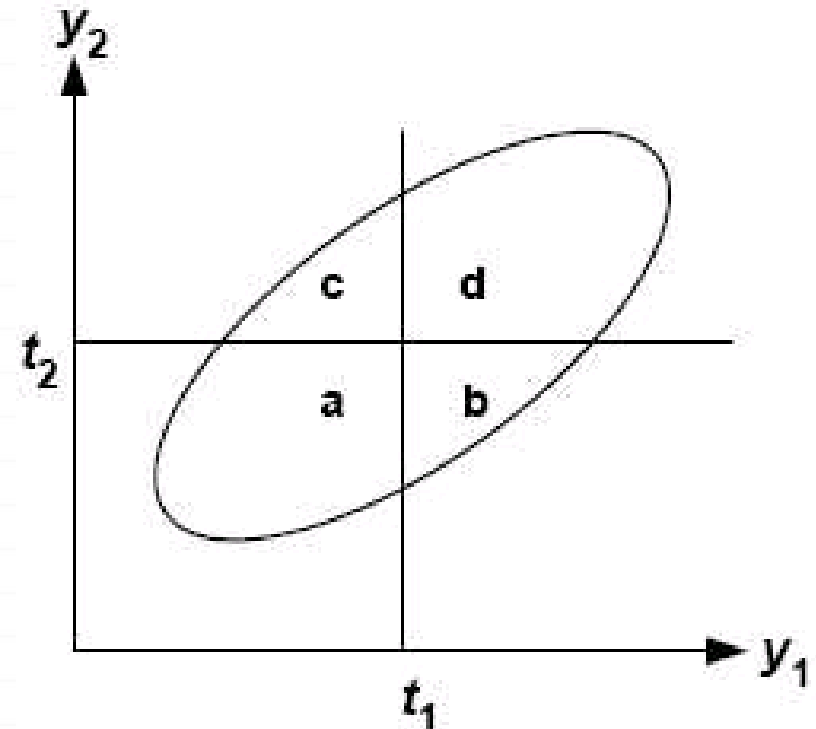
- o With continuous data, we analyse sample correlation matrix
- o To compute correlations for binary data, we refer to *underlying* “quantitative response tendencies” (McDonald, 1999)
 - o These underlying variables are assumed normally distributed
 - o They are connected to the observed responses through a threshold process:

$$\begin{cases} 1 & \text{if } y^* > \tau \\ 0 & \text{if } y^* \leq \tau \end{cases}$$



Tetrachoric correlation

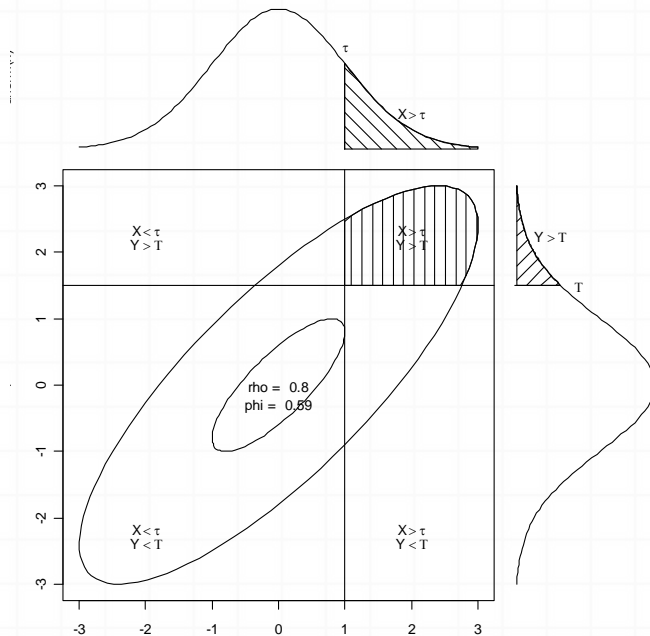
	-	+	y_1
-	a	b	$a + b$
+	c	d	$c + d$
y_2	$a + c$	$b + d$	1



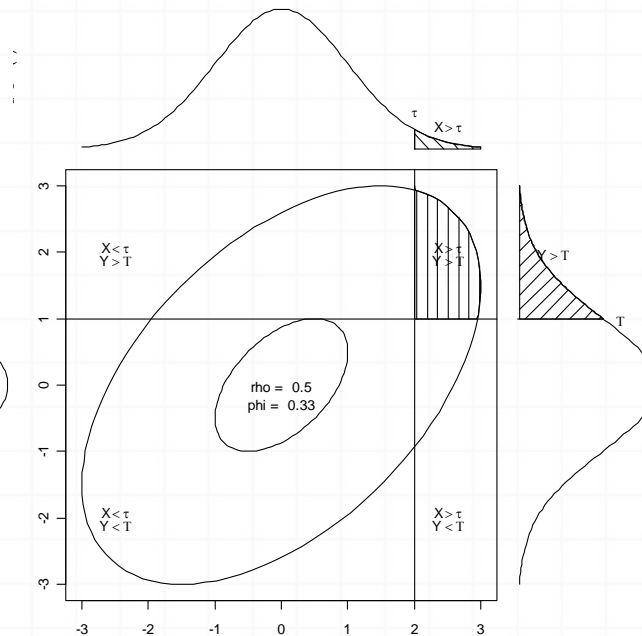
- Tetrachoric correlations can be computed from 2x2 proportions table based on underlying bivariate normal distribution
- Assumes that two normally distributed variables have been dichotomised using a threshold process

Examples of tetrachoric correlations

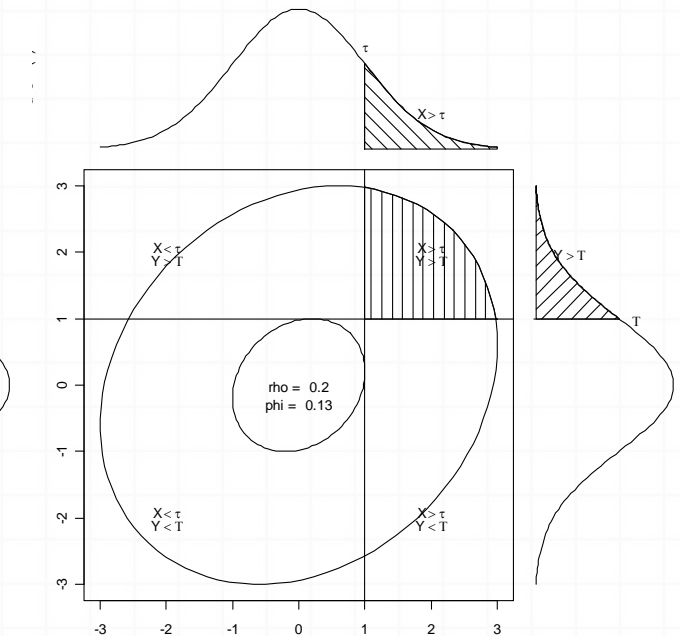
0.8



0.5



0.2



Correlations between ordinal items

- With ordinal data, we have *polychoric correlations*
- Polychorics are used as a convenient *estimation device*, however, for some samples the assumption of multivariate normality might be too strong

CFA extension to categorical variables

- o Mplus provides straightforward extension
- o The only modification to syntax needed is to declare variables as categorical
 - CATEGORICAL ARE i1-i10;
- o Mplus takes care of other things
- o Muthén and Asparauhov (2002) describe an estimation method that considers *underlying response tendencies* (based on tetrachoric or polychoric correlations) to fit a factorial structure

CFA with categorical variables

– setting the scale

1. Continuous variables have scales of their own – categorical variables do not
2. The scale of indicators cannot be passed to their errors (there is nothing to pass)
3. So the scale for errors needs to be set
 1. Mplus sets the error variances to 1 (which is NOT printed anywhere in the output)

Estimation of CFA models with categorical variables

- Default estimator depends on type of analysis and measurement level of observed variables.
- For **categorical** variables default is WLSMV
 - The proper name is Diagonally Weighted Least Squares with mean- and variance-corrected standard errors
 - Makes *no distributional assumptions*
- ULSMV can also be used
- Both are so-called limited information estimators (as opposed to FIML)
- ML can also be used, but it is VERY heavy for more than 3 dimensions, and impossible beyond 4

Example - Inductive reasoning test

- o Fragment of a paper & pencil test assessing aptitude for finding patterns and rules and applying them
- o Consists of ‘passages’ describing different problems (“situations”) – we will consider 5 here
 - o There are 3 problems to solve about each “situation”
- o We analyse data from the test’s first trial, $n=451$

Inductive reasoning test - EFA

TITLE: EFA of Inductive reasoning test

Situations A,B,C,D,E contain 3 questions each

DATA:

FILE IS IndReason.dat; **!individual data**

VARIABLE:

NAMES ARE ID a1-a3 b1-b3 c1-c3 d1-d3 e1-e3;

USEVARIABLES ARE a1-a3 b1-b3 c1-c3 d1-d3 e1-e3;

CATEGORICAL ARE ALL;

ANALYSIS:

TYPE IS EFA 1 5;

ROTATION=CF-VARIMAX (OB); **!we will rotate obliquely**

OUTPUT: RES;

PLOT: TYPE=PLOT3;

Inductive reasoning test – new outputs

- Now Mplus prints out categories proportions and counts

UNIVARIATE PROPORTIONS AND COUNTS FOR CATEGORICAL
VARIABLES

A1

Category 1	0.120	54.000
Category 2	0.880	397.000

A2

Category 1	0.279	126.000
Category 2	0.721	325.000

A3

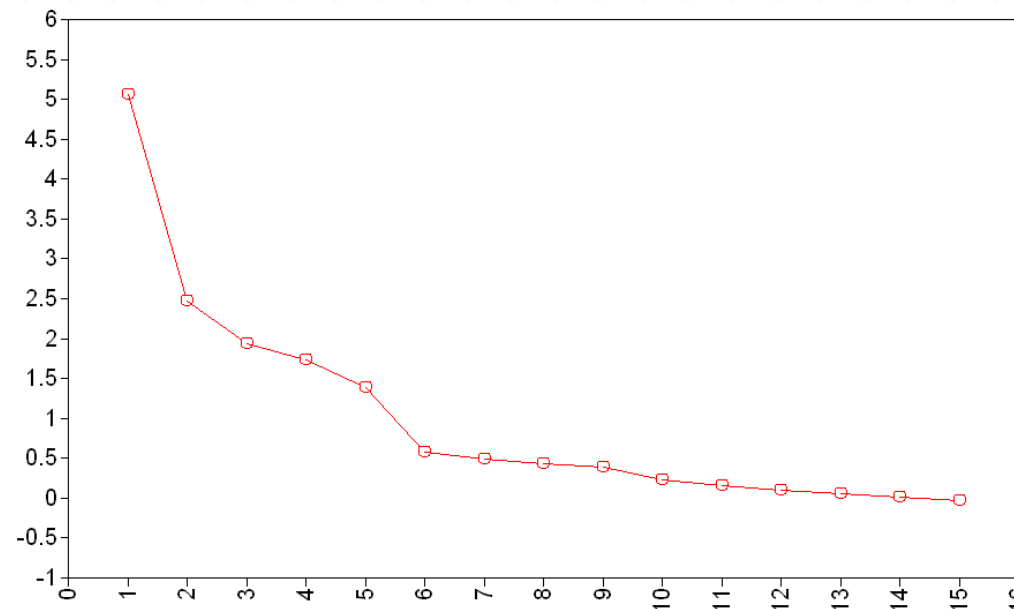
Category 1	0.621	280.000
Category 2	0.379	171.000

Inductive reasoning test

How many factors?

	1 factor	2 factors	3 factors	4 factors	5 factors
Chi square	1139.295	715.886	453.095	209.517	40.631
df	90	76	63	51	40
CFI	.775	.863	.917	.966	1
RMSEA	.161	.137	.117	.083	.006

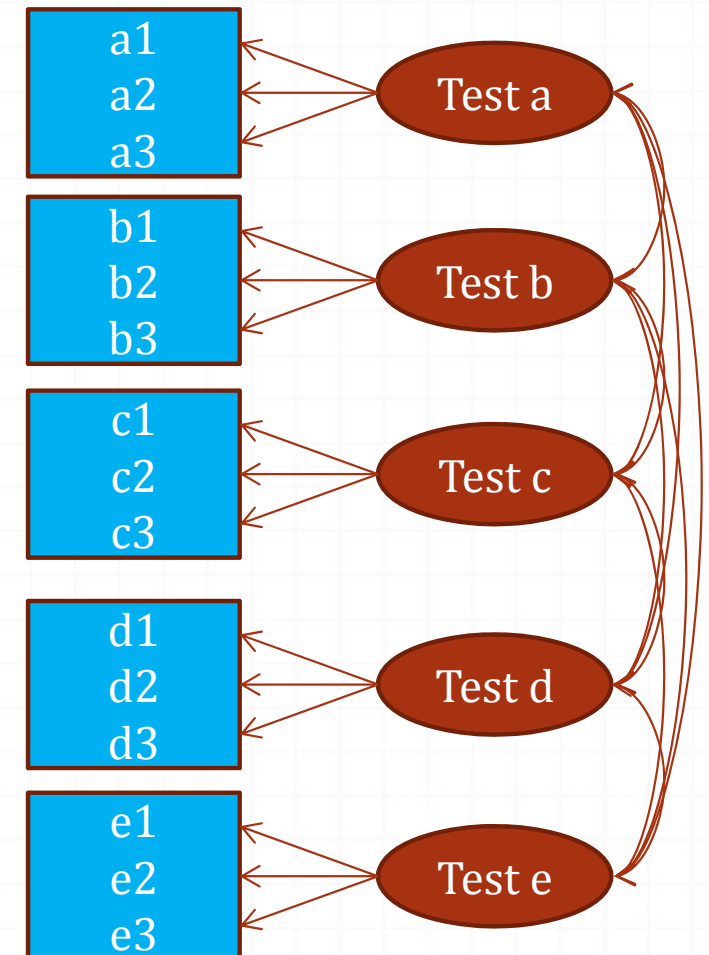
○ Scree plot



Inductive reasoning test

Rotated loadings

	a	b	d	c	e
A1	0.822	0.184	-0.024	0.094	0.047
A2	1.019	-0.038	-0.005	-0.066	-0.002
A3	0.640	0.006	0.127	0.120	-0.034
B1	0.017	0.911	-0.011	0.112	0.045
B2	0.078	0.800	0.072	-0.107	-0.025
B3	0.001	0.601	0.061	0.076	0.068
C1	-0.003	0.043	-0.017	0.801	-0.041
C2	0.026	0.044	-0.001	0.761	0.005
C3	-0.013	-0.008	0.091	0.719	0.081
D1	-0.024	0.002	0.893	0.088	-0.027
D2	0.026	-0.045	0.854	-0.083	0.106
D3	0.028	0.103	0.978	0.042	0.030
E1	-0.062	0.051	0.080	-0.001	0.876
E2	-0.044	0.144	0.007	-0.073	0.911
E3	0.107	-0.069	0.027	0.087	0.980



Inductive reasoning test

Factor correlations

	1	2	3	4	5
1	1.000				
2	0.301	1.000			
3	0.136	0.289	1.000		
4	0.186	0.315	0.215	1.000	
5	0.066	0.290	0.348	0.106	1.000

Inductive reasoning test

EFA model summary

- Standard errors are around 0.05 as they should be; residuals are very small
- Are there really 5 factors? Does each “situation” require a distinct fundamental ability to read and interpret it?
- Or, questions within each “situation” share common variance – *method variance*
 - If the examinee understood the “situation”, all questions relating to it are more likely to be answered correctly (and vice versa)
- This leads to local dependencies of items within “situations” (*correlated uniquenesses*):
 - Common variance in the questions is explained by the overall factor, and unique variance by “situations”

Time for practical #2

CFA and EFA with continuous data in Mplus

Special issues in CFA

When your data is not as simple as textbook examples

Inductive reasoning test - correlated uniquenesses

DATA: FILE IS IndReason.dat; **!individual data**

VARIABLE:

NAMES ARE ID a1-a3 b1-b3 c1-c3 d1-d3 e1-e3;

USEVARIABLES ARE a1-a3 b1-b3 c1-c3 d1-d3 e1-e3;

CATEGORICAL ARE ALL;

MODEL:

REASON BY a1-a3* b1-b3 c1-c3 d1-d3 e1-e3; **!common factor**

REASON@1;

!correlated unique factors related to situations

a1 WITH a2-a3*; a2 WITH a3*;

b1 WITH b2-b3*; b2 WITH b3*;

c1 WITH c2-c3*; c2 WITH c3*;

d1 WITH d2-d3*; d2 WITH d3*;

e1 WITH e2-e3*; e2 WITH e3*;

OUTPUT: RES; MOD; **!request residuals and modification indices**

Inductive reasoning test

Model fit

Chi-Square Test of Model Fit

Value	94.025*
Degrees of Freedom	75
P-Value	0.0679
CFI	0.996
RMSEA	0.024

o Standard errors and residuals are ok

Inductive reasoning test - thresholds

o Now we get estimates of **thresholds** – new output compared to CFA with continuous variables

Thresholds

A1\$1	-1.176	0.077	-15.369	0.000
A2\$1	-0.585	0.063	-9.305	0.000
A3\$1	0.308	0.060	5.125	0.000
...				
D1\$1	0.372	0.061	6.154	0.000
D2\$1	-0.003	0.059	-0.047	0.962
D3\$1	0.625	0.063	9.854	0.000
E1\$1	1.033	0.072	14.345	0.000
E2\$1	0.875	0.068	12.870	0.000
E3\$1	1.112	0.074	14.948	0.000

Inductive reasoning test– standardised factor loadings

REASON	BY	Estimate	S.E.	Est./S.E.	P-Value
A1		0.506	0.080	6.307	0.000
A2		0.236	0.085	2.787	0.005
A3		0.361	0.079	4.586	0.000
B1		0.663	0.087	7.601	0.000
B2		0.510	0.086	5.919	0.000
B3		0.523	0.087	6.019	0.000
C1		0.287	0.084	3.407	0.001
C2		0.350	0.081	4.311	0.000
C3		0.403	0.081	4.995	0.000
D1		0.481	0.082	5.848	0.000
D2		0.426	0.082	5.217	0.000
D3		0.665	0.084	7.930	0.000
E1		0.487	0.100	4.851	0.000
E2		0.475	0.086	5.549	0.000
E3		0.531	0.095	5.600	0.000

Interpreting correlated uniquenesses

- o Normal output will give covariances between residuals
 - o This is useful for evaluating how much residual variance is shared between items from the same “situation”
 - o To evaluate correlations between residuals, one has to examine standardized output
 - o Let’s take item B1
 - o Factor loading .663 (R-square is .439, which means 43.9% of variance is explained by the common “problem solving” factor)
 - o Remaining residual variance is .561; out of which .415 is shared with B2, and .293 is shared with B3. So the “situation” explains roughly as much variance as the common factor.

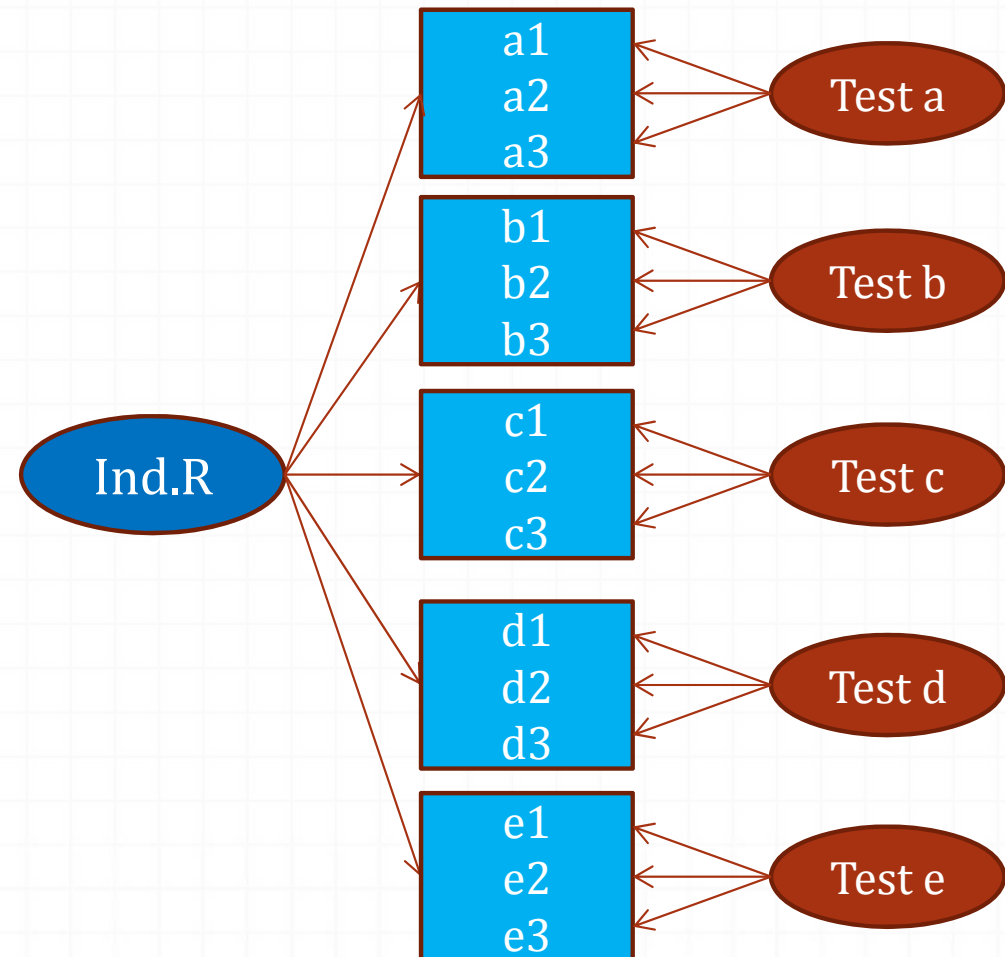
Correlated uniquenesses - issues

- o Estimation of trait scores rests on the assumption of *local independence*
 - o correlated residuals violate this assumption
- o Correlated residuals **always** mean presence of additional factors
 - o This is not understood by all researchers
 - o Correlation between 2 residuals is equivalent to a doublet factor
 - o Which can be modelled
 - o Remember that loadings in doublet factors are not identified?
 - o Mutual correlations between 3 or more residuals may mean a common factor underlying them

Inductive reasoning test

Bifactor model

- This is an alternative to a model with correlated residuals
 - If shared residual variance can be explained by common factors
- A good case for the problem with passages in ability tests
- In a bifactor model, each item loads on 2 factors
 - common factor
 - specific factor



Inductive reasoning test

Bifactor model - syntax

MODEL:

!common factor

REASON BY a1-a3* b1-b3 c1-c3 d1-d3 e1-e3;

REASON@1;

!specific factors

a BY a1-a3*;

b BY b1-b3*;

c BY c1-c3*;

d BY d1-d3*;

e BY e1-e3*;

a-e@1;

!common uncorrelated with specifics, and specifics are uncorrelated with each other

REASON WITH a-e@0;

a-e WITH a-e@0;

Inductive reasoning test

Bifactor model - results

- Fit is the same as for the model with correlated errors (Why?)

Chi-Square 94.025, df=75

- Factor loadings are the same as in model with correlated errors
- Now we get loadings on specific factors, for example

B	BY			
B1	0.688	0.122	5.626	0.000
B2	0.603	0.107	5.652	0.000
B3	0.425	0.102	4.178	0.000

One more alternative for correlated errors - parcels

- Testlets can be also treated as item parcels
- Generally, indicators that have correlated errors are replaced by one indicator that is calculated as their sum
- In our Inductive reasoning test this will constitute count of successes for the whole passage
 - Still a categorical variable!
- Then CFA will proceed with parcels instead of original indicators, and should be free of correlated error problems

Inductive Reasoning – item parcels

- First we compute variables **a, b, c, d** and **e** by using the **DEFINE** command

DEFINE:

```
a=a1+a2+a3;
```

```
b=b1+b2+b3;
```

```
c=c1+c2+c3;
```

```
d=d1+d2+d3;
```

```
e=e1+e2+e3;
```

- Then we use the new variables to test a factor model

MODEL:

```
IndReasoning BY a* b c d e;
```

```
IndReasoning@1;
```

Inductive Reasoning – item parcels model test

Chi-Square Test of Model Fit

Value	10.158*
Degrees of Freedom	5
P-Value	0.0709
RMSEA	0.048
CFI	0.961

o And here are the factor loadings

REASON BY

A	0.379	0.064	5.935	0.000
B	0.612	0.078	7.856	0.000
C	0.363	0.071	5.115	0.000
D	0.510	0.073	6.946	0.000
E	0.474	0.078	6.062	0.000

Another common issue – Nuisance factors

- o Many method factors have been described
 - o Quite often, people agree with items as presented (*acquiescence bias*)
 - o In EFA, 2 factors are found where only 1 should exist
 - o For instance, items assessing Optimism split into 2 clusters – indicating optimism and pessimism
 - o However, optimism and pessimism should be opposite ends of the same factor
- o There are ways of modelling such bias. I will show one example

Random intercept model

- o Recall the standard common factor model (i – item, j – respondent)

$$y_{ij} = m_i + l_i * F_j + e_{ij}$$

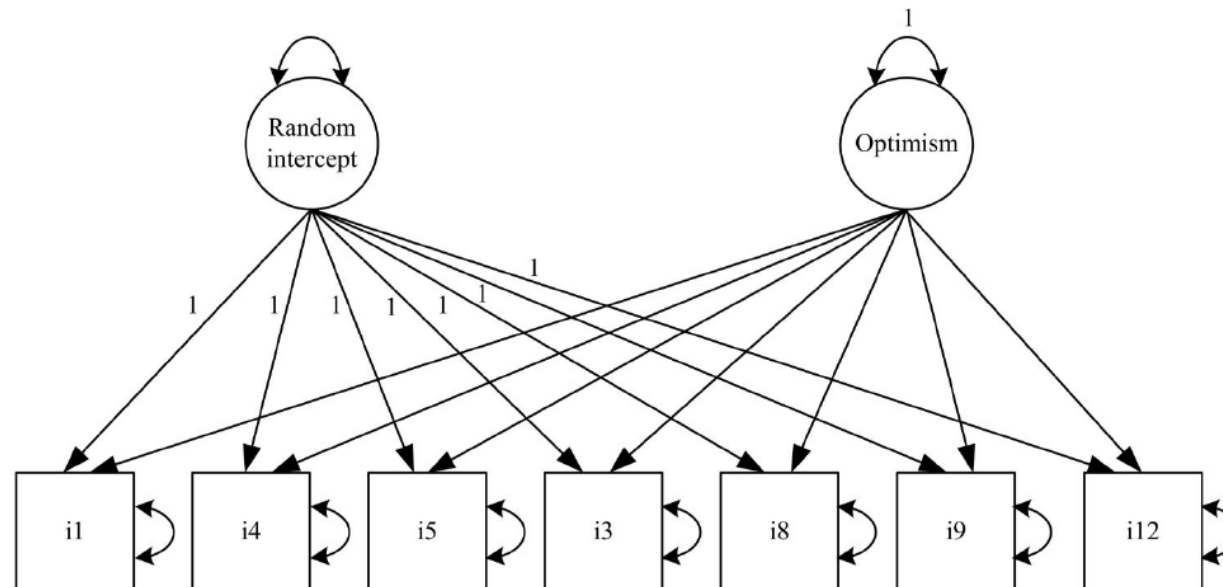
- o The individual tendency to agree (or disagree) with items as presented is incorporated in the model by breaking down the item intercept into a fixed and a random part:

$$y_{ij} = (m_i + \mathbf{RI}_j) + l_i * F_j + e_{ij}$$

- o The fixed part of the intercept varies from item to item
- o The random part is common to all items, but varies from respondent to respondent
 - o If the random part is above zero, the level of agreement with all items is higher than average
 - o If the random part is below zero, the level of agreement with all items is lower than average

Random intercept structural model

- Random intercept is a latent variable that has equal loadings on all items but varies across participants



Reference: Maydeu-Olivares & Coffman (2006). Random intercept factor item analysis. *Psychological Methods*, 11, 344-362.

Syntax for the random intercept model

MODEL:

FACTOR by i1-i20*;

FACTOR@1;

RI BY i1-i20@1; !random intercept has all loadings
equal 1

RI*; !its variance is estimated

FACTOR WITH RI@0;

OUTPUT: RESIDUALS; MODINDICES;

Multi-group CFA

Means and covariance structure

Purpose of multi-group CFA

- Confirmatory approach with multiple groups can be used to test for **any combinations** of the following
 - Measurement parameters (measurement invariance)
 - Equality of Intercepts
 - Equality of Factor loadings
 - Equality of Residual variances
 - Structural parameters (population heterogeneity)
 - Equality of Latent means
 - Equality of Latent variances/covariances
- One of the most attractive features is that **more than 2 groups** can be tested

Defaults for multi-group setup

- o The measurement part of the model **is** assumed invariant if not specified otherwise
 - o Intercepts, factor loadings
 - o (except error variances for continuous indicators) – this is NOT consistent with **strict factorial invariance**
- o The structural part of the model **is not** assumed invariant
 - o Factor means, variances, covariances and regression coefficients

Example – Inductive Reasoning test

- Here we will work with the Inductive reasoning test again
 - Testlets are treated through item parcels
 - Parcel score is assumed categorical (number of successes for the whole passage, ranging from 0 to 3)
- We will analyse data for 2 groups – native English speakers and non-native speakers to see if there are any differences between the groups

Syntax for multi-group analysis

- o Setup for the strictly invariant measurement model

VARIABLE: *<all commands as before>*

GROUPING IS nat_eng (1=native, 2=non-native);

ANALYSIS: PARAM=THETA;

MODEL: REASON BY a b c d e; **!common measurement model**

MODEL non-native: **!group-specific**

a@1 b@1 c@1 d@1 e@1; **!constrain residuals to be the same**

OUTPUT: MODINDICES (ALL 3.84);

- o If we examine the output, it will become obvious which parameters Mplus constrains to be equal across groups

Examining the multi-group output

- The fit of the strict measurement invariance model

Chi-Square Test of Model Fit

Value	41.614*
Degrees of Freedom	28
P-Value	0.0471
RMSEA	0.047 (0.005 0.075)
CFI	0.893

- Examining the modification indices:

REASON BY D MI=8.115

- To free the loading, insert this command to the MODEL section:

MODEL non-native: REASON BY d*;

- Now the model fits: chi-square 28.170 (df=27, p=0.402)

- Loading is **2.271** for native group and **0.373** (n/s) for non-native

Structural parameters

- o Measurement part
 - o Factor loadings, thresholds and residual variances are the same across groups (apart from **Test d** factor loading)
 - o The test is not measurement invariant
- o Factor means and variances
 - o Native speakers **mean= 0** (fixed), **var=0.124**
 - o Non-native speakers **mean = -0.278**, **var = 0.173**
- o Looks like the non-native speakers group might be different in terms of both their mean and variance

Testing for equality of means and variances

- Imposing parameter constraints (one by one)

MODEL:

REASON BY a b c d e; !overall part

REASON (1); !this will enforce equality of variances

MODEL:

REASON BY a b c d e; !overall part

[REASON@0]; !this will enforce equality of means

- Then looking if the fit is worse than fit of the basic model
 - Use DIFFTEST command for WLSMV estimator*
 - The variances are not significantly different
 - The means are different (highly significant chi-square difference)

*Refer to Mplus manual

Continuation of practical #2

CFA and EFA with categorical data and
Multi-group CFA in Mplus