



Exploratory Factor Analysis II

Mengqian Lu

Exploratory Factor Analysis

- ▶ Exploratory Factor Analysis is used to determine the number of **latent variables** that are needed to explain the correlations among a set of observed variables.
- ▶ EFA is to discover the factor structure of a measure and to examine its internal reliability.
- ▶ Recommended when researchers have no hypotheses about the nature of the underlying factor structure of their measure.
- ▶ Three main decision points of EFA: (1) number of factors; (2) extraction method; (3) rotation method.



Example: Ability and Intelligence Tests

- ▶ General model for one individual:

$$x_1 = \mu_1 + \lambda_{11}f_1 + \dots + \lambda_{1q}f_q + u_1$$

...

$$x_p = \mu_p + \lambda_{p1}f_1 + \dots + \lambda_{pq}f_q + u_p$$

$$\text{cov}(x_i) = \Sigma = \Lambda\Lambda^T + \Psi$$

To be determined from x:

1. q : no. of common factors
2. Factor loadings Λ
3. Ψ , the $\text{cov}(\mathbf{u})$
4. Factor scores f

- ▶ Matrix notation for one individual:

$$x = \mu + \Lambda f + u$$

- ▶ Matrix notation for n individuals:

$$x_i = \mu + \Lambda f_i + u_i \quad (i = 1, \dots, n)$$



Factor Models (Cont'd)

► Generalized k -factor model:

Or

$$X_{ji} = \mu_j + \sum_{\ell=1}^k \lambda_{j\ell} F_{\ell i} + U_{ji}, \quad 1 \leq j \leq p, 1 \leq i \leq n$$

$$\mathbf{X}_i = \boldsymbol{\mu} + \boldsymbol{\Lambda}_{(p \times k)} \mathbf{F}_{i(k \times 1)} + \mathbf{U}_i, \quad 1 \leq i \leq n$$

where X_{ji} is the j^{th} random variable in a p -dimensional random vector \mathbf{X}_i , the i^{th} subject.

- $\boldsymbol{\mu} \in \mathbb{R}^p$ constant
- $\boldsymbol{\Lambda}_{(p \times k)}$ is factor loading matrix, explaining the relation between the factors in \mathbf{F}_i and the observed \mathbf{X}_i
- Errors \mathbf{U}_i with $E(\mathbf{U}_i) = 0$; $\text{Cov}(\mathbf{U}_i) = \text{diag}(\psi_1, \dots, \psi_p) = \Psi$
- $\text{Cov}(\mathbf{F}_i, \mathbf{U}_i) = \mathbf{0}$, or $\text{Cov}(u_j, f_s) = 0$ for all j, s Assumptions

Convention: factors are scaled, otherwise $\boldsymbol{\Lambda}$ and $\boldsymbol{\mu}$ are not well determined.

- $\mathbf{F}_{i(k \times 1)}$ is a k -variate random vector with $E(\mathbf{F}_i) = \mathbf{0}$; $\text{Cov}(\mathbf{F}_i) = \mathbf{I}_k$

FA: Covariance matrix

$$\boldsymbol{x} = \boldsymbol{\mu} + \boldsymbol{\Lambda}f + \boldsymbol{u} \Leftrightarrow \text{cov}(\boldsymbol{x}) = \boldsymbol{\Sigma} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^T + \boldsymbol{\Psi}$$

- ▶ Factor model is essentially a particular structure imposed on covariance matrix

Still remember the decomposition of covariance matrix in PCA?

$$\begin{aligned} \text{Cov}(\boldsymbol{x}) &= \frac{1}{n} \sum_i \boldsymbol{x}_i \boldsymbol{x}_i^T = \boldsymbol{C}_x \\ \boldsymbol{C}_x &= \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^T \end{aligned}$$

- Rotations: PCs
- Scores: describe relation between individuals and PCs

FA:

$$\text{cov}(\boldsymbol{x}) = \boldsymbol{\Sigma} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^T + \boldsymbol{\Psi}$$



FA: Covariance matrix

$$\mathbf{x} = \mu + \Lambda \mathbf{f} + \mathbf{u} \Leftrightarrow \text{cov}(\mathbf{x}) = \Sigma = \Lambda \Lambda^T + \Psi$$

- ▶ Factor model is essentially a particular structure imposed on covariance matrix
- ▶ Instead of spectral decomposition in PCA, we look for a way to split up the variances:

$$\text{var}(x_j) = \sigma_j^2 = \text{var}\left(\sum_{k=1}^q \lambda_{jk}^2 F_k\right) + \text{var}(u_j) = \sum_{k=1}^q \lambda_{jk}^2 + \psi_j$$

$\sum_{k=1}^q \lambda_{jk}^2$	Describe the “Communality”, variance due to common factors
ψ_j	Describe the “uniqueness”, variance specific to X_j

- ▶ “Heywood case” in FA, when $\psi_j \leq 0$ (Fabrigar et al.)
-

Factor Analysis in R: `factanal()`

data “ability.cov” in R

general	a non-verbal measure of general intelligence using Cattell's culture-fair test
picture	a picture-completion test
blocks	block design
maze	mazes
reading	reading comprehension
vocab	vocabulary

Ref: Bartholomew, D. J. (1987) and Bartholomew, D. J. and Knott, M. (1990)

In R: `factanal()` perform **maximum-likelihood factor analysis** on a covariance matrix or data matrix



Number of factors

- ▶ MLE approach for estimation provides test:

H_q : q -factor model is sufficient; **vs.** H_u : Σ is unconstrained

- ▶ Practical strategy:
 - Start with a small value of q (e.g. $q=1$) and check hypothesis test; increase q by 1 at a time until some H_q is not rejected

R codes for Intelligence Tests Example

```
ability.FA = factanal(factors = 1, covmat=ability.cov)
update(ability.FA, factors=2)
update(ability.FA, factors=2, rotation="promax")
```




```
> ability.FA
```

```
Call:
```

```
factanal(factors = 1, covmat = ability.cov)
```

```
Uniquenesses:
```

```
general picture blocks maze reading vocab  
  0.535  0.853  0.748  0.910 0.232  0.280
```

```
Loadings:
```

```
          Factor1  
general 0.682  
picture 0.384  
blocks  0.502  
maze    0.300  
reading 0.877  
vocab   0.849
```

```
          Factor1  
SS loadings 2.443  
Proportion Var 0.407
```

```
Test of the hypothesis that 1 factor is sufficient.
```

```
The chi square statistic is 75.18 on 9 degrees of freedom.
```

```
The p-value is 1.46e-12
```



```
> update(ability.FA, factors=2)
```

Call:

```
factanal(factors = 2, covmat = ability.cov)
```

Uniquenesses:

	general	picture	blocks	maze	reading	vocab
	0.455	0.589	0.218	0.769	0.052	0.334

Loadings:

	Factor1	Factor2
general	0.499	0.543
picture	0.156	0.622
blocks	0.206	0.860
maze	0.109	0.468
reading	0.956	0.182
vocab	0.785	0.225

	Factor1	Factor2
SS loadings	1.858	1.724
Proportion Var	0.310	0.287
Cumulative Var	0.310	0.597

Hypothesis can not be rejected; for simplicity, we thus use two factors

Test of the hypothesis that 2 factors are sufficient.

The chi square statistic is 6.11 on 4 degrees of freedom.

The p-value is 0.191

Scale invariance of factor analysis

- ▶ Suppose $y_j = c_j x_j$ or in matrix: $y = Cx$, C is a diagonal matrix, e.g. units conversion, we already have $\Sigma_x = \Lambda_x \Lambda_x^T + \Psi_x$
- ▶ Thus, $Cov(y) = C \Sigma_x C^T = C(\Lambda_x \Lambda_x^T + \Psi_x) C^T$ K-factor model holds for x

$$= (C \Lambda_x)(C \Lambda_x)^T + C \Psi_x C^T = \Lambda_y \Lambda_y^T + \Psi_y$$
i.e. loadings and uniquenesses are the same if expressed in new units

- ▶ In many applications, the search for loadings Λ and for specific variance Ψ is done by the decomposition of the correlation matrix of X rather than the covariance matrix Σ . This corresponds to a factor analysis of a linear combination of X , i.e., $Y = D^{-1}(X - \mu)$, $D = \text{diag}(\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp})$. As you find the Λ_y and Ψ_y , then we have $\Lambda_x = D \Lambda_y$ and $\Psi_x = D \Psi_y D$

Scale invariance of factor analysis (Cont'd)

- ▶ Using *cov* or *cor* gives basically the same result
- ▶ Common practice: use correlation matrix or scale input data (this is done in `factanal()`)



Rotational invariance of factor analysis

► The factor loadings are not unique!

- Suppose that $G^T G = I$ (G : orthogonal matrix), transform f to $f^* = G^T f$, Λ to $\Lambda^* = \Lambda G$, for X^*

- This yields the same model:

$$X^* = \Lambda^* f^* + u = (\Lambda G)(G^T f) + u = \Lambda f + u = X$$

$$\Sigma^* = \Lambda^* \Lambda^{*T} + \Psi = (\Lambda G)(\Lambda G)^T + \Psi = \Lambda \Lambda^T + \Psi = \Sigma$$

- The rotated model is equivalent for explaining the covariance matrix
- Usage: use rotation that makes interpretation of loadings easy or look for rotation that make more practice sense
- Most popular rotation: **Varimax** – each factor has few large and many small loadings



> update(ability.FA, factors=2) Default rotation = 'varimax'

Call:

```
factanal(factors = 2, covmat = ability.cov)
```

Uniquenesses:

	general	picture	blocks	maze	reading	vocab
	0.455	0.589	0.218	0.769	0.052	0.334

Loadings:

	Factor1	Factor2
general	0.499	0.543
picture	0.156	0.622
blocks	0.206	0.860
maze	0.109	0.468
reading	0.956	0.182
vocab	0.785	0.225

Spatial reasoning

Not clear

Interpretation of factors is generally debatable

Verbal intelligence

	Factor1	Factor2
SS loadings	1.858	1.724
Proportion Var	0.310	0.287
Cumulative Var	0.310	0.597

Test of the hypothesis that 2 factors are sufficient.

The chi square statistic is 6.11 on 4 degrees of freedom.

The p-value is 0.191

```
> update(ability.FA, factors=2, rotation="promax")
```

Call:

```
factanal(factors = 2, covmat = ability.cov, rotation = "promax")
```

Uniquenesses:

general	picture	blocks	maze	reading	vocab
0.455	0.589	0.218	0.769	0.052	0.334

Loadings:

	Factor1	Factor2
general	0.364	0.470
picture		0.671
blocks		0.932
maze		0.508
reading	1.023	
vocab	0.811	

Spatial reasoning

Verbal intelligence

	Factor1	Factor2
SS loadings	1.853	1.807
Proportion Var	0.309	0.301
Cumulative Var	0.309	0.610

BETTER? But wait,
orthogonality is not ensured

Factor Correlations:

	Factor1	Factor2
Factor1	1.000	0.557
Factor2	0.557	1.000

Test of the hypothesis that 2 factors are sufficient.

The chi square statistic is 6.11 on 4 degrees of freedom.

The p-value is 0.191

How to find unique solution?

- **Options 1:** Have $\mathbf{M} = \mathbf{\Lambda}\mathbf{\Psi}^{-1}\mathbf{\Lambda}$ be diagonal, with diagonal elements in descending order of magnitude \rightarrow ordered orthogonal factors with descending contributions (same logic as PCA)
- **Potential problems:**
 1. Variables may have substantial loadings on >1 factor
 2. From the 2nd factor, they often turn to be bipolar, i.e., a mixture of positive and negative loadings \rightarrow hard to interpret
- ❖ **Possible solution:** rotate the loadings instead, but there has been debate over this, in my opinion, there is nothing wrong to introduce domain knowledge



How to find unique solution? (Cont'd)

▪ **Options 2: Factor rotation** such that (1) each variable is highly loaded on at most one factor; and/or (2) all factor loadings are either large positive or near zero, with as few intermediate values as possible, thus we split original variables into disjoint sets, each set is associated with a single factor.

→ **Simple Structure** [Thurstone, 1931]

1. *Each row of the factor loading matrix should contain at least one zero.*
2. *Each column of the loading matrix should contain at least k zeros.*
3. *Every pair of columns of the loading matrix should contain several variables whose loadings vanish in one column but not in the other.*
4. *If the number of factors is four or more, every pair of columns should contain a large number of variables with zero loadings in both columns.*
5. *Conversely, for every pair of columns of the loading matrix only a small number of variables should have non-zero loadings in both columns.*

(IAMA by Everitt and Hothorn, page 145)



```
> update(ability.FA, factors=2, rotation="promax")
```

Call:

```
factanal(factors = 2, covmat = ability.cov, rotation = "promax")
```

Uniquenesses:

general	picture	blocks	maze	reading	vocab
0.455	0.589	0.218	0.769	0.052	0.334

Loadings:

	Factor1	Factor2
general	0.364	0.470
picture		0.671
blocks		0.932
maze		0.508
reading	1.023	
vocab	0.811	

	Factor1	Factor2
SS loadings	1.853	1.80
Proportion Var	0.309	0.30
Cumulative Var	0.309	0.61

Factor Correlations:

	Factor1	Factor2
Factor1	1.000	0.557
Factor2	0.557	1.000

Test of the hypothesis that 2
The chi square statistic is 6
The p-value is 0.191

Simple Structure

- 1. Each row of the factor loading matrix should contain at least one zero.*
- 2. Each column of the loading matrix should contain at least k zeros.*
- 3. Every pair of columns of the loading matrix should contain several variables whose loadings vanish in one column but not in the other.*
- 4. If the number of factors is four or more, every pair of columns should contain a large number of variables with zero loadings in both columns.*
- 5. Conversely, for every pair of columns of the loading matrix only a small number of variables should have non-zero loadings in both columns.*

Factor rotation

- ▶ *varimax, promax and more ...*
- ▶ *Still remember? → we may have to abandon this restriction now*

Convention: factors are scaled, otherwise Λ and μ are not well determined.

$F_{i(k \times 1)}$ is a k -variate random vector with $E(F_i) = 0$; $\text{Cov}(F_i) = I_k$

Two main types of rotation:

- ▶ **Orthogonal rotation:** restrict the rotated factors to be uncorrelated
 - Post-multiply the original matrix of loadings by an **orthogonal matrix**
 - After rotation, $\text{Cov}(F_i) = I_k$ still holds
- ▶ **Oblique rotation:** allow correlated factors
 - Post-multiply the original matrix of loadings by an **matrix (no longer constrained to be orthogonal)**
 - After rotation, $\text{Cov}(F_i)$ has unit elements on its diagonal, but no restrictions on the off-diagonal elements



So which one to use, Orthogonal or Oblique?

- ▶ If objective is to “best fit” data, oblique is better; if objective is generalizability, then choose orthogonal to avoid: “what accounts for the relationship between the factors?”
- An orthogonal rotation is simple, an oblique rotation introduce a matrix of factor correlation to consider.
- *Varimax* is orthogonal rotation; *promax* is oblique rotation.
- Factor rotation is still controversial, as it is criticized as – possible allowing one to impose on the data a preconceived structure.



Factor scores

- ▶ Scores are assumed to be random variables
 - ▶ Scores are predicted values for each subject (e.g. individuals who took the intelligence tests)
 - ▶ They represent the original data in a reduced dimension for further analysis or modeling (PCs from PCA also can be used to develop a model for prediction using original data information)
 - ▶ They are likely to be more reliable than the observed variable values – “noises” eliminated
 - ▶ The factor score is a “pure” measure of a latent variable, while an observed value may be unclear due to the “noises”, e.g. nervousness due to unnoticed in-class quiz
-



Factor scores – estimation

- ▶ Scores are assumed to be random variables

- ▶ Two methods to do estimation:
 1. “Bartlett” – Assume f as fix (Maximum-likelihood estimate) – Bartlett's weighted least-squares scores, e.g. `factanal(~v1+v2+v3+v4+v5+v6, factors = 3, scores = "Bartlett")$scores`
 2. “Thompson ” – Assume f as random (Bayesian estimate), e.g. `factanal(~v1+v2+v3+v4+v5+v6, factors = 3, scores = "regression")$scores`
 - No big difference in practice



Example

```
> v1 = c(1,1,1,1,1,1,1,1,1,1,1,3,3,3,3,3,4,5,6)
> v2 = c(1,2,1,1,1,1,2,1,2,1,3,4,3,3,3,4,6,5)
> v3 = c(3,3,3,3,3,1,1,1,1,1,1,1,1,1,1,5,4,6)
> v4 = c(3,3,4,3,3,1,1,2,1,1,1,1,2,1,1,5,6,4)
> v5 = c(1,1,1,1,1,3,3,3,3,3,1,1,1,1,1,6,4,5)
> v6 = c(1,1,1,2,1,3,3,3,4,3,1,1,1,2,1,6,5,4)
> m1 = cbind(v1,v2,v3,v4,v5,v6)
> factanal(m1, factors = 3,
+          scores = "Bartlett")$scores
```

	Factor1	Factor2	Factor3
[1,]	-0.9039949	-0.9308984	0.9475392
[2,]	-0.8685952	-0.9328721	0.9352330
[3,]	-0.9082818	-0.9320093	0.9616422
[4,]	-1.0021975	-0.2529689	0.8178552
[5,]	-0.9039949	-0.9308984	0.9475392
[6,]	-0.7452711	0.7273960	-0.7884733
[7,]	-0.7098714	0.7254223	-0.8007795
[8,]	-0.7495580	0.7262851	-0.7743704
[9,]	-0.8080740	1.4033517	-0.9304636
[10,]	-0.7452711	0.7273960	-0.7884733
[11,]	0.9272282	-0.9307506	-0.8371538
[12,]	0.9626279	-0.9327243	-0.8494600
[13,]	0.9229413	-0.9318615	-0.8230509
[14,]	0.8290256	-0.2528211	-0.9668378
[15,]	0.9272282	-0.9307506	-0.8371538
[16,]	0.4224366	2.0453079	1.2864761
[17,]	1.4713902	1.2947716	0.5451562
[18,]	1.8822320	0.3086244	1.9547752

FA vs. PCA

- ▶ PCA targets explaining **variances**, FA targets explaining **correlations**
 - ▶ PCA is exploratory and without assumptions;
FA is based on statistical model with assumptions
 - ▶ First few PCs will be same regardless of q ;
First few factors of FA depend on q
 - ▶ FA: Orthogonal rotation of factor loadings are equivalent.
This does not hold in PCA
 - ▶ More mathematically:
PCA: $x = \mu + \Gamma_1 z_1 + \Gamma_2 z_2 = \mu + \Gamma_1 z_1 + e$ Assume we only keep the PCs in Γ_1
FA: $x = \mu + \Lambda f + u$
Cov(u) is diagonal by assumption; Cov(e) is not
 - ▶ **! Both PCA and FA only useful if input data is correlated !**
-



FA vs. PCA – Covariance Matrix

As a consequence, $\text{cov}(\mathbf{x}_i) = \Sigma = \Lambda\Lambda^T + \Psi$

- Ψ is diagonal matrix
- We need to estimate this **covariance matrix**

Covariance matrix again, relation to PCA?

	FA	PCA
Covariance	$\text{Cov}(\mathbf{X}_i) = \Sigma = \Lambda\Lambda^T + \Psi$	$\text{Cov}(\mathbf{X}_i) = \mathbf{U}\Lambda\mathbf{U}^T = \mathbf{U}\Lambda^{1/2}(\mathbf{U}\Lambda^{1/2})^T$
Representation	$\mathbf{X} - \mu\mathbf{1}_n^T = \Lambda_{(p \times k)} \mathbf{F}_{(k \times n)} + \mathbf{U}$	$\tilde{\mathbf{X}} \approx \mathbf{U}\mathbf{D}\mathbf{V}^T = \mathbf{U}_{(p \times r)} \mathbf{Z}_{(r \times n)}$
Interpretation	<ul style="list-style-type: none">▪ Factor loading λ describe the relation between an observed X and k hidden factors F_j▪ \mathbf{X} is a linear combination of F_j with weights λ	<ul style="list-style-type: none">▪ PC loading u in the PC direction \mathbf{u}, describe the relation between a PC: Z and the original variable X_j▪ Z is a linear combination of X_j with weights u_j

