

Multivariate Analysis

Factor Analysis

Shaobo Jin

Department of Mathematics

Intended Learning Outcome

Through this chapter, you should be able to

- ① Perform factor analysis,
- ② Explain indeterminacy in factor analysis,
- ③ Obtain factor scores.

Motivation: Latent Variable Modelling

	--	-	-/+	+	++
1. I am a 'worrier'	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
2. I make friends easily	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
3. I have a vivid imagination	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
4. I trust others	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
5. I complete tasks successfully	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
6. I get angry easily	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
7. I really enjoy large parties and gatherings	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
8. I think art is important	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

The Model

The factor analysis model for a $p \times 1$ random vector \mathbf{X} is

$$X_i = \mu_i + \ell_{i1}F_1 + \ell_{i2}F_2 + \cdots + \ell_{im}F_m + e_i, \quad i = 1, 2, \dots, p,$$

where F 's are the **common factors**, ℓ 's are the **(factor) loadings**, e_i is the **error** (or **specific factor**). In the model, we should have $m < p$.

In matrix notation,

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e},$$

where \mathbf{L} is the **factor loading matrix** of size $p \times m$.

Scale Indeterminacy

The factor analysis model is indeterminate, since we do not observe \mathbf{F} .

- ① $\mathbb{E}(\mathbf{F})$ is not unique: for any $m \times 1$ \mathbf{a} ,

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e} = \underbrace{\boldsymbol{\mu} + \mathbf{L}\mathbf{a}}_{\boldsymbol{\mu}^*} + \underbrace{\mathbf{L}(\mathbf{F} - \mathbf{a})}_{\mathbf{F}^*} + \mathbf{e}.$$

- ② $\text{cov}(\mathbf{F})$ is not unique: for any $m \times m$ invertible matrix \mathbf{T} ,

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e} = \boldsymbol{\mu} + \underbrace{(\mathbf{L}\mathbf{T})}_{\mathbf{L}^*} \underbrace{(\mathbf{T}^{-1}\mathbf{F})}_{\mathbf{F}^*} + \mathbf{e}.$$

Assumptions

Our assumptions are

$$\begin{aligned}\mathbb{E}(\mathbf{F}) &= \mathbf{0}, \\ \text{cov}(\mathbf{F}) &= \mathbf{I}, \\ \mathbb{E}(\mathbf{e}) &= \mathbf{0}, \\ \text{cov}(\mathbf{e}) &= \mathbf{\Psi} \text{ (diagonal)}, \\ \text{cov}(\mathbf{F}, \mathbf{e}) &= \mathbf{0}.\end{aligned}$$

The model

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{LF} + \mathbf{e}$$

satisfying the above assumptions is an [orthogonal factor model](#).

Model Implied Covariance Matrix

The orthogonal factor model implies

- ① an expectation for \mathbf{X} : $\mathbb{E}(\mathbf{X}) = \boldsymbol{\mu} + \mathbf{L}\mathbb{E}(\mathbf{F}) + \mathbb{E}(\mathbf{e}) = \boldsymbol{\mu}$.
- ② a covariance matrix for \mathbf{X} :

$$\begin{aligned}\text{cov}(\mathbf{X}) &= \text{cov}(\mathbf{LF} + \mathbf{e}) \\ &= \text{cov}(\mathbf{LF}) + \text{cov}(\mathbf{e}) \\ &= \mathbf{LL}^T + \boldsymbol{\Psi},\end{aligned}$$

which is called a **model-implied covariance matrix**.

This means that

$$\begin{array}{lcl}\sigma_{ii} & = & \sum_{j=1}^m \ell_{ij}^2 \quad + \quad \psi_i \\ \text{variance} & & \text{communality} \quad \quad \text{uniqueness (or specific variance)}\end{array}.$$

We often denote $h_i = \sum_{j=1}^m \ell_{ij}^2$.

Existence of Decomposition

The orthogonal factor analysis model is basically:

Whether we can decompose Σ into $\mathbf{LL}^T + \Psi$, where $\Psi > 0$ is diagonal?

Valid Decomposition

$$\begin{bmatrix} 19 & 30 & 2 \\ 30 & 57 & 5 \\ 2 & 5 & 38 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ -1 \end{bmatrix} \begin{bmatrix} 4 & 7 & -1 \end{bmatrix} + \begin{bmatrix} 2 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Invalid Decomposition

$$\begin{bmatrix} 1 & 0.9 & 0.7 \\ 0.9 & 1 & 0.4 \\ 0.7 & 0.4 & 1 \end{bmatrix} = \begin{bmatrix} 1.2549 \\ 0.7173 \\ 0.5577 \end{bmatrix} \begin{bmatrix} 1.2549 \\ 0.7173 \\ 0.5577 \end{bmatrix}^T + \begin{bmatrix} -0.5748 & 0 & 0 \\ 0 & 0.4852 & 0 \\ 0 & 0 & 0.6891 \end{bmatrix}.$$

Indeterminacy

Even though we let $\text{cov}(\mathbf{F}) = \mathbf{I}$, the model

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e}$$

is still indeterminate. Let \mathbf{T} be any $m \times m$ orthogonal matrix. Then,

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e} = \boldsymbol{\mu} + \underbrace{(\mathbf{L}\mathbf{T})}_{\mathbf{L}^*} \underbrace{(\mathbf{T}^T \mathbf{F})}_{\mathbf{F}^*} + \mathbf{e},$$

$$\text{cov}(\mathbf{X}) = \mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi} = (\mathbf{L}\mathbf{T})(\mathbf{L}\mathbf{T})^T + \boldsymbol{\Psi},$$

$$\text{cov}(\mathbf{F}^*) = \mathbf{T}^T \text{cov}(\mathbf{F}) \mathbf{T} = \text{cov}(\mathbf{F}) = \mathbf{I}.$$

Hence, \mathbf{L} is determined only up to an orthogonal matrix \mathbf{T} .

Scale Invariant

Consider the model

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e}.$$

Suppose that we make a linear transformation $\mathbf{Z} = \mathbf{V}^{-1/2}(\mathbf{X} - \boldsymbol{\mu})$.
Then the covariance matrix of \mathbf{Z} satisfies

$$\begin{aligned}\text{cov}(\mathbf{Z}) &= \mathbf{V}^{-1/2} \text{cov}(\mathbf{X}) \mathbf{V}^{-1/2} \\ &= \left(\mathbf{V}^{-1/2} \mathbf{L}\right) \left(\mathbf{V}^{-1/2} \mathbf{L}\right)^T + \mathbf{V}^{-1/2} \boldsymbol{\Psi} \mathbf{V}^{-1/2},\end{aligned}$$

which correspond to the model

$$\mathbf{Z} = \mathbf{V}^{-1/2} \mathbf{L}\mathbf{F} + \mathbf{V}^{-1/2} \mathbf{e}.$$

Popular Estimation Methods

Suppose that we have observed $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$. We can compute the sample covariance matrix \mathbf{S} or sample correlation matrix \mathbf{R} .

Popular methods of parameter estimation for a factor analysis model are

- ① principal component (or principal factor),
- ② maximum likelihood,
- ③ least squares.

Spectral Decomposition

The spectral decomposition implies that the covariance matrix Σ satisfies

$$\Sigma = \sum_{i=1}^p \lambda_i \mathbf{e}_i \mathbf{e}_i^T = \underbrace{\begin{bmatrix} \sqrt{\lambda_1} \mathbf{e}_1 & \cdots & \sqrt{\lambda_p} \mathbf{e}_p \end{bmatrix}}_{\mathbf{L} \text{ is } p \times p} \begin{bmatrix} \sqrt{\lambda_1} \mathbf{e}_1^T \\ \vdots \\ \sqrt{\lambda_p} \mathbf{e}_p^T \end{bmatrix} + \mathbf{0}.$$

If the last $p - m$ eigenvalues are small, we can ignore them and approximate Σ by

$$\Sigma \approx \underbrace{\begin{bmatrix} \sqrt{\lambda_1} \mathbf{e}_1 & \cdots & \sqrt{\lambda_m} \mathbf{e}_m \end{bmatrix}}_{\mathbf{L} \text{ is } p \times m} \begin{bmatrix} \sqrt{\lambda_1} \mathbf{e}_1^T \\ \vdots \\ \sqrt{\lambda_m} \mathbf{e}_m^T \end{bmatrix} + \begin{bmatrix} \psi_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \psi_p \end{bmatrix},$$

where $\psi_i = \sigma_{ii} - \sum_{j=1}^m \ell_{ij}^2$.

Principal Component Method

Suppose that we have computed the sample covariance matrix \mathbf{S} or the sample correlation matrix \mathbf{R} .

- The estimated factor loading matrix is

$$\hat{\mathbf{L}} = \begin{bmatrix} \sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 & \sqrt{\hat{\lambda}_2} \hat{\mathbf{e}}_2 & \cdots & \sqrt{\hat{\lambda}_m} \hat{\mathbf{e}}_m \end{bmatrix}.$$

- The estimated specific variances are

$$\hat{\psi}_i = s_{ii} - \sum_{j=1}^m \hat{\ell}_{ij}^2.$$

- The diagonal elements of \mathbf{S} are reproduced by $\hat{\mathbf{L}}\hat{\mathbf{L}}^T + \hat{\mathbf{\Psi}}$, but not the off-diagonal elements.

One advantage of the principal component method is that adding more factors does not change the loadings of early factors.

Determine Number of Factors

To determine the number of factors m ,

- 1 subject knowledge,
- 2 using residual matrix: the entries in $\mathbf{S} - \hat{\mathbf{L}}\hat{\mathbf{L}}^T - \hat{\mathbf{\Psi}}$ are small.
- 3 Kaiser criterion: number of eigenvalues of \mathbf{R} greater than 1. Easy but dirty.
- 4 cumulative eigenvalues: the proportion of total sample variance due to j th factor is

$$\frac{\hat{\lambda}_j}{\sum_{i=1}^p s_{ii}}.$$

Choose m such that a suitable proportion of the total sample variance has been explained.

Principal Factor/Axis Method

Consider the sample correlation matrix \mathbf{R} . Suppose that the specific variance ψ_i is ψ_i^* . Then, the communality satisfies

$$h_i^2 = 1 - \psi_i^*.$$

The reduced sample correlation matrix is

$$\mathbf{R}_r = \begin{bmatrix} h_1^{*2} & r_{12} & \cdots & r_{1p} \\ r_{12} & h_2^{*2} & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1p} & r_{2p} & \cdots & h_p^{*2} \end{bmatrix},$$

It should satisfy $\mathbf{R}_r \approx \mathbf{L}_r \mathbf{L}_r^T$, and our loading matrix is \mathbf{L}_r .

Principal Factor/Axis Method: Algorithm

Algorithm 1: Pseudo code of the principal axis method

```
1 Assign initial values  $\psi_i^*$  for all  $i$ , e.g., from eigen decomposition ;
2 Let the number of factors be  $m$  ;
3 while Convergence not met do
4   Compute the reduced correlation matrix  $\mathbf{R}_r$  using  $\mathbf{R}$  and  $\psi_i^*$ 's ;
5   Apply spectral decomposition to  $\mathbf{R}_r$  ;
6   The loading matrix is  $\mathbf{L}_r = \begin{bmatrix} \sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 & \sqrt{\hat{\lambda}_2} \hat{\mathbf{e}}_2 & \cdots & \sqrt{\hat{\lambda}_m} \hat{\mathbf{e}}_m \end{bmatrix}$  ;
7   Update  $\psi_i^*$  using  $\mathbf{L}_r \mathbf{L}_r^T$  ;
8 end
```

Maximum Likelihood

Now we also assume that

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{F} + \mathbf{e} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}).$$

Then, the log-likelihood function is

$$\begin{aligned}\log L(\boldsymbol{\theta}) &= \text{constant} - \frac{n}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \\ &= \text{constant} - \frac{n}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^n \text{tr} \left[\boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \right],\end{aligned}$$

where $\boldsymbol{\theta}$ is the vector of unknown parameters.

Maximum Likelihood Estimator

The MLE maximizes the log-likelihood function

$$\log L(\boldsymbol{\theta}) = \text{constant} - \frac{n}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^n \text{tr} \left[\boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^T \right].$$

- ① The MLE of $\boldsymbol{\mu}$ is $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}}$.
- ② The MLEs of \mathbf{L} and $\boldsymbol{\Psi}$ must be obtained by numerical methods from

$$\begin{aligned} \log L(\boldsymbol{\theta}) &= \text{constant} - \frac{n}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{i=1}^n \text{tr} \left[\boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^T \right] \\ &= \text{constant} - \frac{n}{2} \log |\boldsymbol{\Sigma}| - \frac{n}{2} \text{tr} (\boldsymbol{\Sigma}^{-1} \mathbf{S}_n). \end{aligned}$$

- ① Since the model is not identified, we often use the restriction $\mathbf{L}^T \boldsymbol{\Psi} \mathbf{L}$ is diagonal and the diagonal entries are written in decreasing order.

Large Sample Test for the Number of Common Factors

An advantage of maximum likelihood is that we can test

H_0 : m factors are sufficient to describe the data

H_1 : Σ has no constraints

The likelihood ratio statistic is given by

$$\begin{aligned} -2 \log \Lambda &= -2 \log \frac{\exp \left\{ -\frac{n}{2} \log |\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\Psi}| - \frac{n}{2} \text{tr} \left[\left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\Psi} \right)^{-1} \mathbf{S}_n \right] \right\}}{\exp \left\{ -\frac{n}{2} \log |\mathbf{S}_n| - \frac{n}{2} \text{tr} \left(\mathbf{S}_n^{-1} \mathbf{S}_n \right) \right\}} \\ &= n \left\{ \text{tr} \left[\left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\Psi} \right)^{-1} \mathbf{S}_n \right] - \log \left| \left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\Psi} \right)^{-1} \mathbf{S}_n \right| - p \right\}, \end{aligned}$$

which has an asymptotic χ_s^2 distribution under H_0 with $s = 2^{-1}p(p+1) - [pm + m - 2^{-1}m(m-1)]$.

Bartlett Correction

In the psychometrics literature,

$$F(\mathbf{\Sigma}, \mathbf{S}_n) = \text{tr}(\mathbf{\Sigma}^{-1} \mathbf{S}_n) - \log |\mathbf{\Sigma}^{-1} \mathbf{S}_n| - p$$

is often called a normal-theory [fit function](#).

The likelihood ratio statistic is then

$$\begin{aligned} -2 \log \Lambda &= n \left\{ \text{tr} \left[\left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\mathbf{\Psi}} \right)^{-1} \mathbf{S}_n \right] - \log \left| \left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\mathbf{\Psi}} \right)^{-1} \mathbf{S}_n \right| - p \right\} \\ &= n F \left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\mathbf{\Psi}}, \mathbf{S}_n \right). \end{aligned}$$

The chi-square approximation is improved if we instead consider

$$\left[n - 1 - \frac{1}{6} (2p + 4m + 5) \right] F \left(\hat{\mathbf{L}} \hat{\mathbf{L}}^T + \hat{\mathbf{\Psi}}, \mathbf{S}_n \right).$$

Other Estimation Methods

We can use fit functions other than the fit function derived from the log-likelihood.

- Normal-theory fit function:

$$F(\boldsymbol{\Sigma}, \mathbf{S}_n) = \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{S}_n) - \log |\boldsymbol{\Sigma}^{-1} \mathbf{S}_n| - p$$

- Unweighted least squares:

$$F(\boldsymbol{\Sigma}, \mathbf{S}_n) = \frac{1}{2} \text{tr} \left\{ (\mathbf{S}_n - \boldsymbol{\Sigma})^2 \right\}.$$

- Generalized least squares:

$$F(\boldsymbol{\Sigma}, \mathbf{S}_n) = \frac{1}{2} \text{tr} \left\{ [(\mathbf{S}_n - \boldsymbol{\Sigma}) \mathbf{S}_n^{-1}]^2 \right\}.$$

For most estimators from fit functions, $\hat{\boldsymbol{\Psi}}$ is not guaranteed to be positive definite. It is called a **Heywood case**.

Parallel Analysis

The number of factors can also be determined using a simulation-based [parallel analysis](#).

Algorithm 2: Pseudo code of the principal axis method

- 1 Fit a factor model using the data and obtain eigenvalues of \mathbf{LL}^T ;
 - 2 **while** *Simulation* **do**
 - 3 Generate random data of same dimension as raw data by either resampling the raw data or simulating from (independent) normal data ;
 - 4 Obtain eigenvalues using the random data ;
 - 5 **end**
 - 6 Plot the eigenvalues from the real data and the random data ;
-

Versus PCA

FA	PCA
Dimension reduction	Dimension reduction
Reflective, with causal interpretation	Formative
Model assumes latent causal variables exist.	Model does not assume latent variables exist.
Scale invariant	Not scale invariant
$\Sigma = \mathbf{L}\mathbf{L}^T + \Psi$	Spectral decomposition

Orthogonal Rotation

The model-implied covariance matrix of a factor model is

$$\Sigma = \mathbf{L}\mathbf{L}^T + \Psi.$$

For any $m \times m$ orthogonal matrix \mathbf{T} , the model-implied covariance matrix is unchanged:

$$\Sigma = \mathbf{L}\mathbf{L}^T + \Psi = (\mathbf{L}\mathbf{T})(\mathbf{L}\mathbf{T})^T + \Psi,$$

and the factor model is equivalent to

$$\mathbf{X} = \boldsymbol{\mu} + (\mathbf{L}\mathbf{T})\mathbf{T}^T\mathbf{F} + \mathbf{e}$$

with $\text{cov}(\mathbf{T}^T\mathbf{F}) = \mathbf{T}^T\text{cov}(\mathbf{F})\mathbf{T} = \text{cov}(\mathbf{F}) = \mathbf{I}$.

Orthogonal Rotation

We can multiply \mathbf{L} by an orthogonal matrix \mathbf{T} (equivalently \mathbf{F} by \mathbf{T}^T) to rotate the factor to ease interpretation. Let \mathbf{B} be the rotated loading matrix.

- **Varimax rotation** maximizes

$$\sum_{k=1}^m \left\{ \frac{1}{p} \sum_{i=1}^p b_{ik}^4 - \left(\frac{1}{p} \sum_{i=1}^p b_{ik}^2 \right)^2 \right\},$$

the columnwise variance of the squared factor loadings.

- **Quartimax rotation** minimizes

$$- \sum_{i=1}^p \sum_{k=1}^m b_{ik}^4.$$

Oblique rotation

Suppose that \mathbf{T} is nonsingular, but not orthogonal. Then,

$$\mathbf{\Sigma} = \mathbf{L}\mathbf{L}^T + \mathbf{\Psi} = (\mathbf{L}\mathbf{T})\mathbf{T}^{-1}(\mathbf{T}^{-1})^T(\mathbf{L}\mathbf{T})^T + \mathbf{\Psi},$$

and the factor model is equivalent to

$$\mathbf{X} = \boldsymbol{\mu} + (\mathbf{L}\mathbf{T})\mathbf{T}^{-1}\mathbf{F} + \mathbf{e}.$$

In this model,

$$\text{cov}(\mathbf{T}^{-1}\mathbf{F}) = \mathbf{T}^{-1}\text{cov}(\mathbf{F})(\mathbf{T}^{-1})^T = \mathbf{T}^{-1}(\mathbf{T}^{-1})^T.$$

This is the **oblique rotation** with new factor $\mathbf{T}^{-1}\mathbf{F}$ and new loading matrix $\mathbf{L}\mathbf{T}$. The covariance matrix of the rotated factor is not \mathbf{I} .

Oblique rotation

Let \mathbf{B} be the rotated loading matrix.

- **Promax rotation** first applies the varimax rotation, and then rotate the initial loading matrix towards the varimax rotated matrix.
- **Quartimin rotation** minimizes

$$\sum_{\ell=1}^m \sum_{k \neq \ell}^p \sum_{i=1}^p b_{ik}^2 b_{i\ell}^2.$$

- **Geomin rotation**, given a small $\epsilon > 0$ (e.g., $\epsilon = 0.01$), minimizes

$$\sum_{i=1}^p \left[\prod_{k=1}^m (b_{ik}^2 + \epsilon) \right]^{1/m}.$$

Bartlett Score

Factor scores are the estimated values of the common factors.

For every subject, the Bartlett score minimizes

$$L(\mathbf{f}) = (\mathbf{x} - \boldsymbol{\mu} - \mathbf{L}\mathbf{f})^T \boldsymbol{\Psi}^{-1} (\mathbf{x} - \boldsymbol{\mu} - \mathbf{L}\mathbf{f}).$$

The minimizer is

$$\mathbf{f}_B = (\mathbf{L}^T \boldsymbol{\Psi}^{-1} \mathbf{L})^{-1} \mathbf{L}^T \boldsymbol{\Psi}^{-1} (\mathbf{x} - \boldsymbol{\mu}).$$

Plugging in the estimators,

$$\hat{\mathbf{f}}_B = \left(\hat{\mathbf{L}}^T \hat{\boldsymbol{\Psi}}^{-1} \hat{\mathbf{L}} \right)^{-1} \hat{\mathbf{L}}^T \hat{\boldsymbol{\Psi}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}).$$

Thompson Score

Suppose that

$$\begin{bmatrix} \mathbf{F} \\ \mathbf{X} \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ \boldsymbol{\mu} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Phi} & \boldsymbol{\Phi} \mathbf{L}^T \\ \mathbf{L} \boldsymbol{\Phi} & \boldsymbol{\Sigma} \end{bmatrix} \right).$$

Then,

$$\mathbf{F} \mid \mathbf{X} = \mathbf{x} \sim N \left(\boldsymbol{\Phi} \mathbf{L}^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}), \boldsymbol{\Phi} - \boldsymbol{\Phi} \mathbf{L}^T \boldsymbol{\Sigma}^{-1} \mathbf{L} \boldsymbol{\Phi} \right).$$

The **Thompson score** (or **regression score**) is

$$\begin{aligned} \mathbf{f}_R &= \boldsymbol{\Phi} \mathbf{L}^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= \left(\mathbf{L}^T \boldsymbol{\Psi}^{-1} \mathbf{L} + \boldsymbol{\Phi}^{-1} \right)^{-1} \mathbf{L}^T \boldsymbol{\Psi}^{-1} (\mathbf{x} - \boldsymbol{\mu}). \end{aligned}$$

Plugging in the estimators,

$$\hat{\mathbf{f}}_R = \left(\hat{\mathbf{L}}^T \hat{\boldsymbol{\Psi}}^{-1} \hat{\mathbf{L}} + \hat{\boldsymbol{\Phi}}^{-1} \right)^{-1} \hat{\mathbf{L}}^T \hat{\boldsymbol{\Psi}}^{-1} (\mathbf{x} - \hat{\boldsymbol{\mu}}).$$

Issues of Factor Scores

However,

$$\begin{aligned}\text{cov}(\mathbf{f}_B) &= (\mathbf{L}^T \mathbf{\Psi}^{-1} \mathbf{L})^{-1} \mathbf{L}^T \mathbf{\Psi}^{-1} \mathbf{\Sigma} \mathbf{\Psi}^{-1} \mathbf{L} (\mathbf{L}^T \mathbf{\Psi}^{-1} \mathbf{L})^{-1}, \\ \text{cov}(\mathbf{f}_R) &= \mathbf{\Phi} \mathbf{L}^T \mathbf{\Sigma}^{-1} \mathbf{L} \mathbf{\Phi},\end{aligned}$$

which are not necessary equal to $\mathbf{\Phi}$.

We cannot use the factor scores as if they were observed. If they are used in a regression model, they should be treated as [errors-in-variables](#).

Anderson-Rubin Score

The Anderson-Rubin method minimizes

$$\sum_{j=1}^n (\mathbf{x}_j - \boldsymbol{\mu} - \mathbf{L}\mathbf{f}_j)^T \boldsymbol{\Psi}^{-1} (\mathbf{x}_j - \boldsymbol{\mu} - \mathbf{L}\mathbf{f}_j)$$
$$\text{s.t. } \frac{1}{n-1} \sum_{j=1}^n \mathbf{f}_j \mathbf{f}_j^T = \text{cov}(\mathbf{F}).$$

Hence, the sample covariance matrix of the resulting factor score is equal to the (estimated) covariance matrix of the factors from the model.

Regression with Factor Score

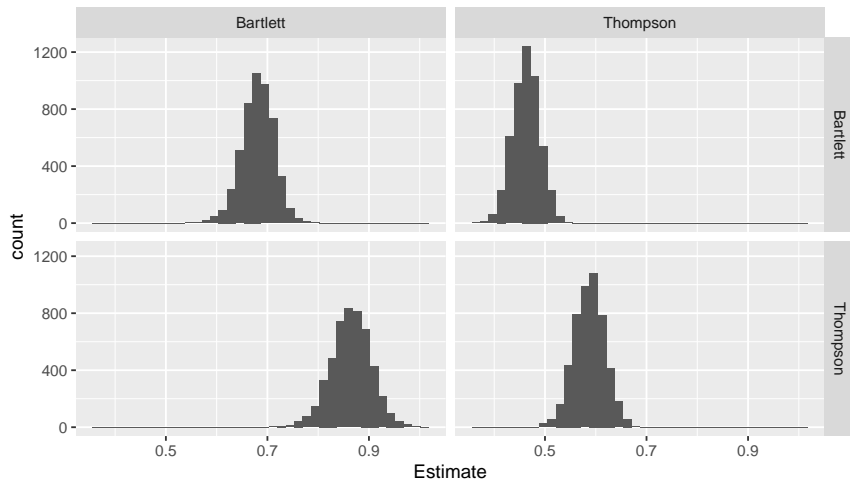
I generate 5000 data sets each with 400 observations from the model

$$\begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.8 \\ 0.7 \\ 0.6 \\ 0.5 \end{bmatrix} F_X + \begin{bmatrix} e_1^X \\ e_2^X \\ e_3^X \\ e_4^X \end{bmatrix}, \quad \begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.8 \\ 0.7 \\ 0.6 \\ 0.5 \end{bmatrix} F_Y + \begin{bmatrix} e_1^Y \\ e_2^Y \\ e_3^Y \\ e_4^Y \end{bmatrix}$$

and

$$F_Y = 0 + 0.75F_X + e.$$

Biased Results



Final Remarks

- The estimation methods that we introduced are mostly designed for continuous \mathbf{X} . If you have a discrete \mathbf{X} , other methods need to be used.
- The factor analysis that we introduced is often called **exploratory factor analysis**, as compared to **confirmatory factor analysis**.
 - Confirmatory factor analysis knows *a priori* which factor loadings should be zero, and no factor rotation is needed.