

## Chapter 2: Factor Analysis<sup>a</sup>

Factor models (FMs) are **statistical models** that attempt to explain the correlation amongst the observed variables  $x_{(1)}^0, \dots, x_{(p)}^0$  via a small number  $k < p$  of **unobservable factors**.

Like PCA, FMs can be used to obtain a low dimensional representation of the observations (where low means  $k$ -dimensional).

FMs are also useful to infer unobservable variables or abstract concepts from proxy variables, such as

- intelligence from test results,
- ‘Big Five’ personality traits from personality surveys,
- salient attributes determining perception of a product.

Throughout this chapter we assume that the observations  $\{x_i^0\}_{i=1}^n$  are  $n$  realizations of an  $\mathbb{R}^p$ -valued random variable  $X^0$ .

---

<sup>a</sup>The main reference for this chapter is [2, Chapter 9].

## The factor model

For  $k < p$  the  $k$ -factor model assumes that the random variable  $X := X^0 - \mathbb{E}[X^0]$  is such that  $X = \Lambda F + U$  where

- $\Lambda = [\lambda_{jl}]$  is a  $p \times k$  matrix of constants, called **loadings matrix**,
- $F$  is an  $\mathbb{R}^k$ -valued random variable, called **factor**, such that  $\mathbb{E}[F] = 0$  and such that  $\text{Var}(F) = \mathbf{I}_k$ ,
- $U$  is an  $\mathbb{R}^p$ -valued random variable such that  $\mathbb{E}[U] = 0$  and such that  $\text{Var}(U) = \Psi$  for some matrix  $\Psi = \text{diag}(\psi_{11}, \dots, \psi_{pp})$ ,
- $\text{Cov}(F, U) = 0$ .

**Remark:** If  $\text{Var}(X)$  is full rank then the FM always holds true if  $k = p$ , in which case  $\Lambda$  and  $F$  can be obtained using population PCA<sup>a</sup>. This is why in this chapter we focus on the case where  $k < p$ .

Under the  $k$ -FM, for all  $j \in \{1, \dots, p\}$  we have for the  $j$ th variable

$$X_j = \sum_{l=1}^k \lambda_{jl} F_l + U_j \Rightarrow \text{Var}(X_j) = \sum_{l=1}^k \lambda_{jl}^2 + \psi_{jj}$$

where

- $h_j^2 := \sum_{l=1}^k \lambda_{jl}^2$  is called **communality**; this is the part of the variance of  $X_j$  that is shared with the other variables through the common factor  $F$ .
- $\psi_{jj} \geq 0$  is called the **specific variance**; this is the part of the variance of  $X_j$  which is not shared with the other variables.

---

<sup>a</sup>To see this recall that  $X = \Gamma Y$  where  $\Gamma \in O(p)$  and where  $Y$  is such that  $\text{Var}(Y) = \mathbf{L}$ , with  $\mathbf{L}$  a diagonal matrix having the eigenvalues of  $\text{Var}(X)$  as non-zero entries (see Chapter 1, page 24). If  $\text{Var}(X)$  is full rank the matrix  $\mathbf{L}$  is invertible and  $X = (\Gamma \mathbf{L}^{1/2})(\mathbf{L}^{-1/2} Y)$ , showing that for  $k = p$  the FM holds with  $\Lambda = \Gamma \mathbf{L}^{1/2}$ ,  $F = \mathbf{L}^{-1/2} Y$  and  $\Psi = 0 \mathbf{I}_p$ .

## Variance decomposition and the fundamental theorem of factor analysis

Let  $\Sigma = \text{Var}(X)$  and note that the  $k$ -factor model implies that  $\Sigma$  can be decomposed as follows:

$$\Sigma = \Lambda \text{Var}(F) \Lambda^\top + \text{Var}(U) = \Lambda \Lambda^\top + \Psi. \quad (2.1)$$

Interestingly, as shown in Theorem 2.1 below, the  $k$ -factor model holds for  $X$  **if and only if** its variance  $\Sigma$  can be decomposed as in (2.1) (under a weak assumption on  $\Sigma$ ).

### Theorem 2.1 (Fundamental theorem of factor analysis)

*Assume that there exists a  $k < p$  such that*

$$\Sigma = \tilde{\Lambda} \tilde{\Lambda}^\top + \tilde{\Psi} \quad (2.2)$$

*for some  $p \times k$  matrix  $\tilde{\Lambda}$  and some  $p \times p$  diagonal matrix  $\tilde{\Psi}$  having strictly positive diagonal entries. Then, the  $k$ -factor model holds for  $X$  with  $\Lambda = \tilde{\Lambda}$  and with  $\Psi = \tilde{\Psi}$ .*

**Remark:** Theorem 2.1 notably assumes that  $\Sigma$  is full rank.

**Remark:** In Theorem 2.1, for a fix matrix  $\tilde{\Psi}$  the decomposition (2.2) is unique up to a rotation of  $\tilde{\Lambda}$ ; that is,

$$\Sigma = \tilde{\Lambda} \tilde{\Lambda}^\top + \tilde{\Psi} = M M^\top + \tilde{\Psi}$$

if and only if  $M = \tilde{\Lambda} G$  for some  $G \in O(k)$  [see 2, Section 9.2.3].

## Proof of Theorem 2.1

To simplify the notation let  $\mathbf{\Lambda} = \tilde{\mathbf{\Lambda}}$  and  $\mathbf{\Psi} = \tilde{\mathbf{\Psi}}$ , so that  $\mathbf{\Sigma} = \mathbf{\Lambda}\mathbf{\Lambda}^\top + \mathbf{\Psi}$ . Remark that under the assumptions of the theorem the matrices  $\mathbf{\Sigma}$  and  $\mathbf{\Psi}$  are invertible.

Let

$$\mathbf{M} = \begin{pmatrix} \mathbf{I}_p & \mathbf{\Lambda} \\ -\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} & \mathbf{I}_k \end{pmatrix}$$

and note that  $\mathbf{M}$  is invertible since

$$\begin{aligned} \det(\mathbf{M}) &= \det(\mathbf{I}_k) \det(\mathbf{I}_p - \mathbf{\Lambda}\mathbf{I}_k^{-1}(-\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1})) \\ &= \det((\mathbf{\Psi} + \mathbf{\Lambda}\mathbf{\Lambda}^\top)\mathbf{\Psi}^{-1}) = \det(\mathbf{\Sigma}) \det(\mathbf{\Psi}^{-1}) > 0. \end{aligned}$$

Let  $Y \sim \mathcal{N}_k(0, \mathbf{I}_k + \mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{\Lambda})$  be independent of  $X$  and let

$$\begin{pmatrix} U \\ F \end{pmatrix} = \mathbf{M}^{-1} \begin{pmatrix} X \\ Y \end{pmatrix}.$$

Then,

$$\mathbf{M} \begin{pmatrix} U \\ F \end{pmatrix} = \begin{pmatrix} U + \mathbf{\Lambda}F \\ -\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} U + F \end{pmatrix} = \begin{pmatrix} X \\ Y \end{pmatrix}$$

from which we obtain that  $X = \mathbf{\Lambda}F + U$ .

To complete the proof it remains to verify that  $\mathbb{E}[F] = 0$ ,  $\mathbb{E}[U] = 0$ ,  $\text{Var}(U) = \mathbf{\Psi}$  and that  $\text{Cov}(F, U) = 0$ .

First, we have  $\mathbb{E}[(U, F)] = \mathbf{M}^{-1} \mathbb{E}[(X, Y)] = (0, 0)$  as required.

To proceed further recall that if  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$ , are four matrices such that  $\mathbf{D}$  and  $\mathbf{E} := \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}$  are invertible matrices then

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{E}^{-1} & -\mathbf{E}^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{E}^{-1} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{E}^{-1}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}. \quad (2.3)$$

Then, applying (2.3) to invert  $\mathbf{M}^{-1}$  we obtain

$$\mathbf{M}^{-1} = \left[ \begin{pmatrix} \mathbf{\Psi} & \mathbf{\Lambda} \\ -\mathbf{\Lambda}^\top & \mathbf{I}_k \end{pmatrix} \begin{pmatrix} \mathbf{\Psi}^{-1} & 0 \\ 0 & \mathbf{I}_k \end{pmatrix} \right]^{-1} = \begin{pmatrix} \mathbf{\Psi} & 0 \\ 0 & \mathbf{I}_k \end{pmatrix} \begin{pmatrix} \mathbf{\Psi} & \mathbf{\Lambda} \\ -\mathbf{\Lambda}^\top & \mathbf{I}_k \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{\Psi} & 0 \\ 0 & \mathbf{I}_k \end{pmatrix} \begin{pmatrix} \mathbf{\Sigma}^{-1} & -\mathbf{\Sigma}^{-1}\mathbf{\Lambda} \\ \mathbf{\Lambda}^\top \mathbf{\Sigma}^{-1} & \mathbf{I}_k - \mathbf{\Lambda}^\top \mathbf{\Sigma}^{-1} \mathbf{\Lambda} \end{pmatrix}.$$

### Proof of Theorem 2.1 (end)

Using the above expression for  $M^{-1}$  we have

$$\begin{aligned}
 \text{Var}((U, F)) &= M^{-1} \text{Var}((X, Y)) (M^{-1})^\top \\
 &= \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & -\Sigma^{-1}\Lambda \\ \Lambda^\top \Sigma^{-1} & I_k - \Lambda^\top \Sigma^{-1} \Lambda \end{pmatrix} \begin{pmatrix} \Sigma & 0 \\ 0 & I_k + \Lambda^\top \Psi^{-1} \Lambda \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & -\Sigma^{-1}\Lambda \\ \Lambda^\top \Sigma^{-1} & I_k - \Lambda^\top \Sigma^{-1} \Lambda \end{pmatrix}^\top \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix} \\
 &= \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix} \begin{pmatrix} I_p & -\Sigma^{-1}\Lambda(I_k + \Lambda^\top \Psi^{-1} \Lambda) \\ \Lambda^\top & (I_k - \Lambda^\top \Sigma^{-1} \Lambda)(I_k + \Lambda^\top \Psi^{-1} \Lambda) \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & \Sigma^{-1}\Lambda \\ -\Lambda^\top \Sigma^{-1} & I_k - \Lambda^\top \Sigma^{-1} \Lambda \end{pmatrix} \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix} \\
 &= \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix} \begin{pmatrix} \text{(i)} & \text{(ii)} \\ \text{(ii)}^\top & \text{(iv)} \end{pmatrix} \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix} = \begin{pmatrix} \Psi & 0 \\ 0 & I_k \end{pmatrix},
 \end{aligned}$$

where

$$\begin{aligned}
 \text{(i)} &= \Sigma^{-1} + \Sigma^{-1}\Lambda(I_k + \Lambda^\top \Psi^{-1} \Lambda)\Lambda^\top \Sigma^{-1} \\
 &= \Sigma^{-1}(\Sigma + \Lambda\Lambda^\top + \Lambda\Lambda^\top \Psi^{-1} \Lambda\Lambda^\top)\Sigma^{-1} \\
 &= \Sigma^{-1}(\Sigma + \Sigma - \Psi + (\Sigma - \Psi)\Psi^{-1}(\Sigma - \Psi))\Sigma^{-1} \\
 &= 2\Sigma^{-1}\Sigma\Sigma^{-1} - \Sigma^{-1}\Psi\Sigma^{-1} + (\Sigma^{-1}\Sigma - \Sigma^{-1}\Psi)\Psi^{-1}(\Sigma\Sigma^{-1} - \Psi\Sigma^{-1}) \\
 &= 2\Sigma^{-1} - \Sigma^{-1}\Psi\Sigma^{-1} + \Psi^{-1} - \Sigma^{-1} - \Sigma^{-1} - \Sigma^{-1}\Psi\Sigma^{-1} = \Psi^{-1} \\
 \text{(ii)} &= \Sigma^{-1}\Lambda - \Sigma^{-1}\Lambda(I_k + \Lambda^\top \Psi^{-1} \Lambda)(I_k - \Lambda^\top \Sigma^{-1} \Lambda) \\
 &= \Sigma^{-1}\Lambda - \Sigma^{-1}\Lambda(I_k + \Lambda^\top \Psi^{-1} \Lambda - \Lambda^\top \Sigma^{-1} \Lambda - \Lambda^\top \Psi^{-1} \Lambda\Lambda^\top \Sigma^{-1} \Lambda) \\
 &= -\Sigma^{-1}\Lambda\Lambda^\top(\Psi^{-1} - \Sigma^{-1} - \Psi^{-1}\Lambda\Lambda^\top \Sigma^{-1})\Lambda \\
 &= -\Sigma^{-1}(\Sigma - \Psi)(\Psi^{-1} - \Sigma^{-1} - \Psi^{-1}(\Sigma - \Psi)\Sigma^{-1})\Lambda \\
 &= -(I_p - \Sigma^{-1}\Psi)(\Psi^{-1} - \Sigma^{-1} - (\Psi^{-1} - \Sigma^{-1}))\Lambda = 0 \\
 \text{(iv)} &= \Lambda^\top \Sigma^{-1} \Lambda + (I_k - \Lambda^\top \Sigma^{-1} \Lambda)(I_k + \Lambda^\top \Psi^{-1} \Lambda)(I_k - \Lambda^\top \Sigma^{-1} \Lambda) \\
 &= \Lambda^\top \Sigma^{-1} \Lambda + (I_k - \Lambda^\top \Sigma^{-1} \Lambda + \Lambda^\top \Psi^{-1} \Lambda - \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top \Psi^{-1} \Lambda)(I_k - \Lambda^\top \Sigma^{-1} \Lambda) \\
 &= \Lambda^\top \Sigma^{-1} \Lambda + I_k - \Lambda^\top \Sigma^{-1} \Lambda + \Lambda^\top \Psi^{-1} \Lambda - \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top \Psi^{-1} \Lambda \\
 &\quad - \Lambda^\top \Sigma^{-1} \Lambda + \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top \Sigma^{-1} \Lambda - \Lambda^\top \Psi^{-1} \Lambda\Lambda^\top \Sigma^{-1} \Lambda + \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top \Psi^{-1} \Lambda\Lambda^\top \Sigma^{-1} \Lambda \\
 &= I_k + \Lambda^\top \Psi^{-1} \Lambda - \Lambda^\top \Sigma^{-1} \Lambda - \Lambda^\top \Psi^{-1} \Lambda\Lambda^\top \Sigma^{-1} \Lambda \\
 &\quad + \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top (\Sigma^{-1} - \Psi^{-1} + \Psi^{-1}\Lambda\Lambda^\top \Sigma^{-1})\Lambda \\
 &= I_k + \Lambda^\top \Psi^{-1} \Lambda - \Lambda^\top \Sigma^{-1} \Lambda - \Lambda^\top \Psi^{-1} (\Sigma - \Psi)\Sigma^{-1} \Lambda \\
 &\quad + \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top (\Sigma^{-1} - \Psi^{-1} + \Psi^{-1}(\Sigma - \Psi)\Sigma^{-1})\Lambda \\
 &= I_k + \Lambda^\top \Psi^{-1} \Lambda - \Lambda^\top \Sigma^{-1} \Lambda - \Lambda^\top (\Psi^{-1} - \Sigma^{-1})\Lambda \\
 &\quad + \Lambda^\top \Sigma^{-1} \Lambda\Lambda^\top (\Sigma^{-1} - \Psi^{-1} + \Psi^{-1} - \Sigma^{-1})\Lambda = I_k.
 \end{aligned}$$

The proof is complete.

## Two important properties of factor analysis

- **Invariance of the factors:** Let  $X' = \mathbf{K}X$  for some diagonal matrix  $\mathbf{K}$ . Then, if  $X = \mathbf{\Lambda}F + U$  we have

$$X' = \mathbf{K}\mathbf{\Lambda}F + \mathbf{K}U, \quad \text{Var}(\mathbf{K}U) = \mathbf{K}\mathbf{\Psi}\mathbf{K}$$

so that the  $k$ -factor model holds for  $X'$  with factor loadings  $\mathbf{K}\mathbf{\Lambda}$ , specific variances  $\mathbf{K}\mathbf{\Psi}\mathbf{K}$  and the same factor  $F$ .

$\Rightarrow$  Unlike PCA, factor analysis (FA) is unaffected by the scaling of the variables (i.e. if the variables are rescaled then the factor  $F$  remains unchanged).

- **Non-Uniqueness of the factor loadings:** For  $\mathbf{G} \in O(k)$  we have

$$X = \mathbf{\Lambda}F + U = (\mathbf{\Lambda}\mathbf{G})(\mathbf{G}^\top F) + U$$

where  $\text{Var}(\mathbf{G}^\top F) = \mathbf{G}^\top \mathbf{G} = \mathbf{I}_k$  and  $\text{Cov}(\mathbf{G}^\top F, U) = 0$ .

$\Rightarrow$  If the FM holds for  $X$  then the loadings matrix  $\mathbf{\Lambda}$  and the factor  $F$  are not unique.

Following the remark made just after Theorem 2.1, for a given  $\mathbf{\Psi}$  the factor loadings are unique up to a rotation. This indeterminacy in the definition of  $\mathbf{\Lambda}$  can be resolved by forcing this matrix to satisfy a constraint, such as [see 2, Section 9.2.3]

$$(a) \mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{\Lambda} \text{ is diagonal or } (b) \mathbf{\Lambda}^\top \mathbf{D}^{-1} \mathbf{\Lambda} \text{ is diagonal,}$$

where in either case the diagonal elements are written e.g. in non-increasing order, and where in (b)  $\mathbf{D}$  is some diagonal matrix.

**Remark:** Conditions (a)-(b) are scale invariant and in (a) we must have  $\psi_{jj} > 0$  for  $j = 1, \dots, p$ .

## Computing the model parameters

Because of the scale invariance of FA the model parameters are often estimated from the correlation matrix  $\mathbf{R}$  of the data (and not from  $\mathbf{S}$ ). Below we follow this common practice.

The goal is then to compute  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$  from  $\mathbf{R}$  in such a way that

$$\mathbf{R} = \mathbf{\Lambda}\mathbf{\Lambda}^\top + \mathbf{\Psi} \quad \text{where } \mathbf{\Lambda} \text{ satisfies (a) or (b).} \quad (2.4)$$

To see if a solution to this problem exists remark that  $\mathbf{\Lambda}$  contains  $kp$  parameters and  $\mathbf{\Psi}$  contains  $p$  parameters. However, under either (a) or (b) there are  $k(k+1)/2 - k = k(k-1)/2$  constraints<sup>a</sup> on  $\mathbf{\Lambda}$ , and thus under (a) or (b) the total number of free parameters in the matrix  $\mathbf{\Lambda}\mathbf{\Lambda}^\top + \mathbf{\Psi}$  is equal to  $n_{\text{fp}} := kp + p - k(k-1)/2$ . On the other hand, the equality  $\mathbf{R} = \mathbf{\Lambda}\mathbf{\Lambda}^\top + \mathbf{\Psi}$  imposes  $n_{\text{cons}} := p(p+1)/2$  constraints on the matrices  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$ .

Let

$$\Delta_{k,p} = n_{\text{cons}} - n_{\text{fp}} = \frac{(p-k)^2}{2} - \frac{(p+k)}{2} \quad (2.5)$$

and consider the following three cases:

- If  $\Delta_{k,p} = 0$  there are as many free parameters as constraints and (2.4) usually has an exact and unique solution.
- If  $\Delta_{k,p} < 0$  there are more free parameters than constraints and (2.4) has infinitely many solutions<sup>b</sup>  $\Rightarrow$  the FM is not well-defined.
- If  $\Delta_{k,p} > 0$  there are more constraints than free parameters and (2.4) has no solution. In this case we look for an approximate solution.

---

<sup>a</sup>This is because both (a) and (b) imposes that all the off diagonal elements of a  $k \times k$  symmetric matrix are equal to zero.

<sup>b</sup>Note that neither (a) nor (b) ensures that  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$  are unique. Instead, under (a) or (b)  $\mathbf{\Lambda}$  is unique given that we know a matrix  $\mathbf{\Psi}$  such that the FM holds (see below).

### Approximate solution to (2.4) when $\Delta_{k,p} > 0$

The scenario  $\Delta_{k,p} > 0$  is the most frequent one in practice, in which case the objective is to compute loadings  $\hat{\Lambda}$  and specific variances  $\hat{\Psi}$  in such a way that  $\mathbf{R} \approx \hat{\Lambda}\hat{\Lambda}^\top + \hat{\Psi}$ , where  $\hat{\Lambda}$  satisfies (a) or (b).

To this aim assume for now that we know a diagonal matrix  $\Psi$  such that the matrix  $\mathbf{R} - \Psi$  is

- (i) positive semi-definite,
- (ii) has rank  $k < p$ .

Let  $\gamma_1 \geq \dots \geq \gamma_p \geq 0$  be the eigenvalues of the matrix  $\mathbf{R} - \Psi$ ,  $\mathbf{B}$  be the corresponding matrix of orthonormal eigenvectors,  $\mathbf{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_p)$  and  $\mathbf{\Gamma}_k = \text{diag}(\gamma_1, \dots, \gamma_k)$ .

Then, using the fact that  $\gamma_{k+1} = \dots = \gamma_p = 0$ , we have

$$\mathbf{R} - \Psi = \mathbf{B}\mathbf{\Gamma}\mathbf{B}^\top = \mathbf{B}_{1:k}\mathbf{\Gamma}_k\mathbf{B}_{1:k}^\top = (\mathbf{B}_{1:k}\mathbf{\Gamma}_k^{1/2})(\mathbf{B}_{1:k}\mathbf{\Gamma}_k^{1/2})^\top$$

showing that if we know  $\Psi$  such that (i) and (ii) hold then the decomposition  $\mathbf{R} = \mathbf{\Lambda}\mathbf{\Lambda}^\top + \Psi$  holds for  $\mathbf{\Lambda} = \mathbf{B}_{1:k}\mathbf{\Gamma}_k^{1/2}$ . Moreover, condition (b) holds for  $\mathbf{D} = \mathbf{I}_p$  since  $\mathbf{\Lambda}^\top\mathbf{\Lambda} = \mathbf{\Gamma}_k$  when  $\mathbf{\Lambda} = \mathbf{B}_{1:k}\mathbf{\Gamma}_k^{1/2}$ .

In practice a diagonal matrix  $\Psi$  such that (ii) holds is of course unknown (and may not exist when  $\Delta_{k,p} > 0$ ), and the idea of the algorithm introduced below is to alternate between computing  $\mathbf{\Lambda}$  as above for a given  $\Psi$  and computing  $\Psi$  for a given matrix  $\mathbf{\Lambda}$ , using  $\Psi = \mathbf{R} - \mathbf{\Lambda}\mathbf{\Lambda}^\top$ .



## The iterated principal factor analysis algorithm

**Input:** A preliminary estimate  $\hat{\Psi}^{(0)}$  of  $\Psi$

**for**  $s \geq 1$  **do**

(i) Compute the eigenvalues  $\gamma_1^{(s)} \geq \dots \geq \gamma_p^{(s)}$  of the **reduced correlation matrix**  $\tilde{\mathbf{R}}^{(s)} := \mathbf{R} - \hat{\Psi}^{(s-1)}$  and a corresponding set  $\{v_j^{(s)}\}_{j=1}^p$  of orthonormal eigenvectors.

(ii) Set  $\hat{\Lambda}^{(s)} = [\hat{\lambda}_{jl}^{(s)}] = \mathbf{B}^{(s)}(\mathbf{\Gamma}^{(s)})^{1/2}$  where

$$\mathbf{B}^{(s)} = [v_1^{(s)} \dots v_k^{(s)}], \quad \mathbf{\Gamma}^{(s)} = \text{diag}(\gamma_1^{(s)}, \dots, \gamma_k^{(s)}).$$

(iii) Let  $\hat{\Psi}^{(s)} = \text{diag}(\psi_{11}^{(s)}, \dots, \psi_{pp}^{(s)})$  where

$$\hat{\psi}_{jj}^{(s)} = \max \left\{ 1 - \sum_{l=1}^k (\hat{\lambda}_{jl}^{(s)})^2, 0 \right\}, \quad \forall j \in \{1, \dots, p\}.$$

**if** Convergence=TRUE **then**

(iv) **return**  $\hat{\Lambda}^{(s)}$  and  $\hat{\Psi}^{(s)}$ .

(v) **break**

**end if**

**end for**

**Remark:** The above definition of  $\hat{\Psi}^{(s)}$  ensures that the equality  $\mathbf{R} = \hat{\Lambda}^{(s)}(\hat{\Lambda}^{(s)})^\top + \hat{\Psi}^{(s)}$  holds at least for the diagonal elements.

**Remark:** A standard choice for  $\hat{\Psi}^{(0)}$  is

$$\hat{\Psi}^{(0)} = \text{diag} \left( 1 - \max_{l \neq 1} |r_{1l}|, \dots, 1 - \max_{l \neq p} |r_{pl}| \right)^a.$$

---

<sup>a</sup>Another possible choice is  $\hat{\Psi}^{(0)} = \text{diag}(1 - \sum_{l \neq 1} |r_{1l}|, \dots, 1 - \sum_{l \neq p} |r_{lp}|)$ . In this case, the initial reduced correlation matrix  $\tilde{\mathbf{R}}^{(1)}$  is diagonal dominant and symmetric, and thus is positive semi-definite.

### Varimax rotation to increase interpretability

As we saw above, the FM is defined up to a rotation of the loadings matrix  $\mathbf{\Lambda}$ , and the aim of the varimax rotation is to find a rotation that makes the estimated FM easy to interpret.

Let  $(\delta_{jl}^{(\mathbf{G})}) = \mathbf{\Lambda}\mathbf{G}$  be the loadings obtained by rotating  $\mathbf{\Lambda}$  using the matrix  $\mathbf{G} \in O(p)$ , and let  $(\tilde{F}_1^{(\mathbf{G})}, \dots, \tilde{F}_k^{(\mathbf{G})}) = \mathbf{G}^\top \mathbf{F}$  be the vector containing the corresponding rotated factors.

Then, the varimax rotation of the loadings  $\mathbf{\Lambda}$  is the matrix  $\mathbf{G}_* \in O(p)$  that maximises the function

$$O(p) \ni \mathbf{G} \mapsto \sum_{l=1}^k \sum_{j=1}^p \left( d_{jl}(\mathbf{G})^2 - \bar{d}_l(\mathbf{G}) \right)^2,$$

where  $d_{jl}(\mathbf{G}) = \delta_{jl}^{(\mathbf{G})}/h_j \in [-1, 1]$  and  $\bar{d}_l(\mathbf{G}) = \frac{1}{p} \sum_{j=1}^p d_{jl}(\mathbf{G})^2$ , recalling that  $h_j^2$  is the  $j$ th diagonal element of  $\mathbf{\Lambda}\mathbf{\Lambda}^\top = (\mathbf{\Lambda}\mathbf{G})(\mathbf{\Lambda}\mathbf{G})^\top$ .

**Remark:**  $d_{jl}(\mathbf{G})$  is the correlation between the factor  $\tilde{F}_l^{(\mathbf{G})}$  and the variable  $X_j$  when  $\psi_{jj} = 0$ .

For all  $l$ , the sample variance of  $d_{1l}(\mathbf{G}_*)^2, \dots, d_{pl}(\mathbf{G}_*)^2$  being large it typically follows that we have  $d_{jl}(\mathbf{G}_*) \approx 0$  for many  $j \in \{1, \dots, p\}$  and  $d_{jl}(\mathbf{G}_*) \approx \pm 1$  otherwise.

Therefore, the loadings matrix  $\mathbf{\Lambda}\mathbf{G}_*$  typically contains a few large loadings (in absolute value) and many near-zero loadings. When this happens the varimax rotation allows to reduce the number of factors on which depends a given component of  $X$ , which facilitates their interpretation.

## Estimation of the factors and dimension reduction

Given an estimate  $(\hat{\Lambda}, \hat{\Psi})$  of  $(\Lambda, \Psi)$  we can estimate the factor  $f_i$  associated to observation  $x_i$  by estimating the coefficient  $\beta \in \mathbb{R}^k$  in the linear regression model

$$X_{ij} = \sum_{l=1}^k \hat{\lambda}_{jl} \beta_l + \epsilon_{ij}, \quad j = 1, \dots, p, \quad \text{Var}(\epsilon_i) = \hat{\Psi}.$$

Using the generalized least squares estimate of  $\beta$ , we obtain the following estimate  $\hat{f}_i$  of the factor  $f_i$ :

$$\hat{f}_i = \mathbf{A}_k^{(\text{FA})} x_i, \quad \mathbf{A}_k^{(\text{FA})} := (\hat{\Lambda}^\top \hat{\Psi}^{-1} \hat{\Lambda})^{-1} \hat{\Lambda}^\top \hat{\Psi}^{-1}.$$

**PCA versus FA:** Recall that the  $q$ -dimensional representation of  $x_i$  obtained with PCA is  $x'_i = \mathbf{A}_{1:q}^\top x_i$ .

- Main differences between PCA and FA:
  - The rows of  $\mathbf{A}_k^{(\text{FA})}$  are not orthogonal, while those of  $\mathbf{A}_{1:q}^\top$  are.
  - $\hat{f}_i$  is model based, while  $x'_i$  is model free.
- Which method to choose for dimension reduction?
  - There is no definitive answer to this question, and quite often both PCA and FA lead to similar results (see the example below).
  - It makes sense to use FA if the FM is a reasonable assumption.

**Remark:** Like PCA, factor analysis can be used as a dimension reduction technique in regression models. This is done by replacing in principal component regression the matrix  $\mathbf{A}_{1:q}^\top$  by the matrix  $\mathbf{A}_k^{(\text{FA})}$ .

### Illustration of FA: The marks dataset (continued)

We consider the marks dataset, already used in the previous chapter on PCA, to illustrate the use of FA. We recall that this dataset contains the marks of  $n = 88$  students in  $p = 5$  subjects.

For this dataset we obtain the following correlation matrix  $\mathbf{R}$ :

variable	Matrix $\mathbf{R}$				
mechanics	1.00	0.55	0.55	0.41	0.39
vectors		1.00	0.61	0.48	0.44
algebra			1.00	0.71	0.66
analysis				1.00	0.61

For  $k > 2$  the value of  $\Delta_{k,p}$  defined in (2.5) is negative, and thus the factor model is well-defined only for  $k \in \{1, 2\}$ .

In what follows we let  $k = 2$  be the number of factors. In this case,  $\Delta_{k,p} > 0$  and an approximate solution to (2.4) is computed using the iterated FA algorithm introduced in this chapter.

### Illustration of FA: The marks dataset (continued)

The matrix  $\hat{\mathbf{\Lambda}} = [\hat{\lambda}_{(1)} \ \hat{\lambda}_{(2)}]$  computed with the iterated FA algorithm is as follows:

variable	$\hat{\lambda}_{(1)}$	$\hat{\lambda}_{(2)}$
mechanics	-0.64	0.34
vectors	-0.71	0.29
algebra	-0.90	-0.09
analysis	-0.77	-0.23
statistics	-0.72	-0.23

while  $\hat{\mathbf{\Psi}} = \text{diag}(0.47, 0.41, 0.19, 0.35, 0.43)$ .

Letting  $\hat{\mathbf{R}} = [\hat{r}_{ij}] = \hat{\mathbf{\Lambda}}\hat{\mathbf{\Lambda}}^\top + \hat{\mathbf{\Psi}}$  and writing  $\mathbf{R}$  as  $\mathbf{R} = [r_{ij}]$ , we observe that  $|r_{ij} - \hat{r}_{ij}| < 0.006$  for all  $i, j \in \{1, \dots, p\}$ . Therefore, the  $k = 2$  factors model fits the data well (i.e. the approximation  $\hat{\mathbf{\Lambda}}\hat{\mathbf{\Lambda}}^\top + \hat{\mathbf{\Psi}} \approx \mathbf{R}$  is good).

To interpret the two factors of the estimated model recall that this latter leads to the approximation  $x_i \approx \hat{\lambda}_{(1)}\hat{f}_{i1} + \hat{\lambda}_{(2)}\hat{f}_{i2}$ .

Since all the components of  $\hat{\lambda}_{(1)}$  are negative it follows that (minus) the first factor can be interpreted as “intelligence” (the larger  $-\hat{f}_{i1}$  the higher the intelligence and thus the higher the marks in all the five subjects).

Recalling that the examination was open book for the first two subjects and closed book for the remaining three subjects, the second factor can be interpreted as the propensity to perform well in one type of examination and poorly in the other.

## FA versus PCA: The marks dataset

PCA with  $q = 2$  principal components leads to the approximation  $x_i \approx a_{(1)}y_{i1} + a_{(2)}y_{i2}$ , where the vectors  $a_{(1)}$  and  $a_{(2)}$  are given in Table 1.1.

From Table 1.1, we observe that  $a_{(1)}$  is very similar to  $\hat{\lambda}_{(1)}$  and that  $a_{(2)}$  is very similar to  $\hat{\lambda}_{(2)}$ . Therefore, if the two approximations  $x_i \approx \hat{\lambda}_{(1)}\hat{f}_{i1} + \hat{\lambda}_{(2)}\hat{f}_{i2}$  and  $x_i \approx a_{(1)}y_{i1} + a_{(2)}y_{i2}$  are good then we must have  $(\hat{f}_{i1}, \hat{f}_{i2}) \approx (y_{i1}, y_{i2})$ , that is, the estimated two dimensional factor  $\hat{f}_i$  must be similar to the  $q = 2$  dimensional reduction  $x'_i$  of  $x_i$  obtained with PCA.

From Figure 2.1 below we remark that for the marks dataset PCA and FA indeed give a very similar two dimensional representation of the data points.

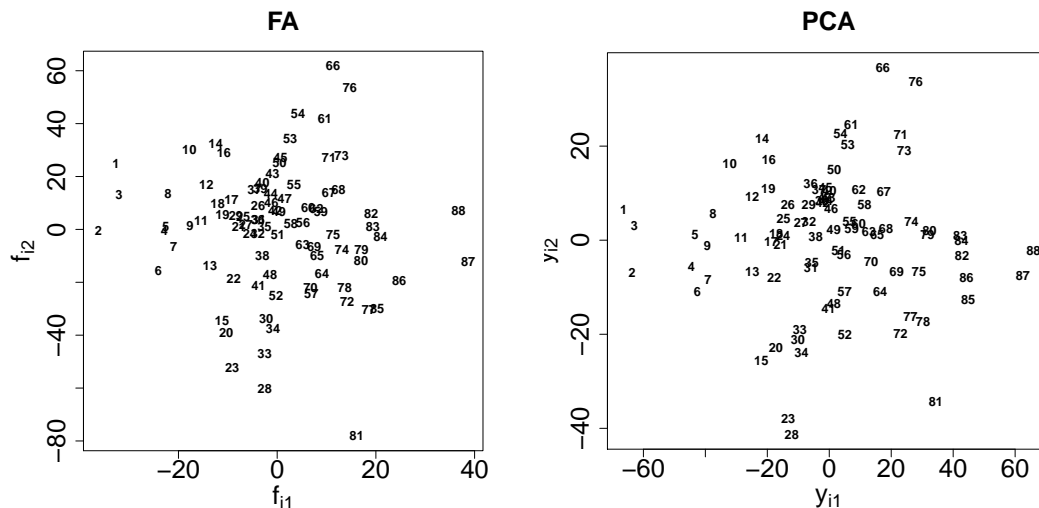


Figure 2.1: Two dimensional approximation of the marks dataset obtained with FA and with PCA.