



COLUMBIA UNIVERSITY  
IN THE CITY OF NEW YORK

STAT 4224/5224

*Bayesian Statistics*

Dobrin Marchev

# **The Factor Analysis Model:**

## **Understanding of Causes**

Factor analysis was invented nearly 100 years ago by psychologist Charles Spearman, who hypothesized that the enormous variety of tests of mental ability measures of mathematical skill, vocabulary, other verbal skills, artistic skills, logical reasoning ability, etc. could all be explained by one underlying "factor" of general intelligence that he called  $g$ . He hypothesized that if  $g$  could be measured and you could select a subpopulation of people with the same score on  $g$ , in that subpopulation you would find no correlations among any tests of mental ability. In other words, he hypothesized that  $g$  was the only factor common to all those measures.

# Factor Analysis

- Factor analysis has a tremendous appeal for the behavior and social sciences.
- In these areas, it is natural to regard multivariate observations on human processes and behavior as manifestations of underlying unobservable “traits”.
- Factor analysis provides a way of explaining the observed variability in behavior in terms of these traits.

# The Factor Analysis Model

- Data:

$\mathbf{X} = (x_1, \dots, x_p)$  is an *observable* random vector which has a  $p$ -variate normal distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ .

- Factors:

$F_1, F_2, \dots, F_m$  are *unobservable/latent* random variables called *the common factors*

- Errors:

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$  are random variables called *errors* or *specific factors*.

- Loadings:

Suppose that there exist constants  $\lambda_{ij}$  (*the loadings*) such that:

$$x_1 = \mu_1 + \lambda_{11}F_1 + \lambda_{12}F_2 + \dots + \lambda_{1m}F_m + \varepsilon_1$$

$$x_2 = \mu_2 + \lambda_{21}F_1 + \lambda_{22}F_2 + \dots + \lambda_{2m}F_m + \varepsilon_2$$

...

$$x_p = \mu_p + \lambda_{p1}F_1 + \lambda_{p2}F_2 + \dots + \lambda_{pm}F_m + \varepsilon_p$$

# Factor Analysis Model in Matrix Notation

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

where

$\mathbf{X}$  is  $p \times 1$ ,  $\mathbf{L}$  is  $p \times m$ ,  $\mathbf{F}$  is  $m \times 1$ , and  $\boldsymbol{\varepsilon}$  is  $p \times 1$

Assume:  $\text{cov}(\mathbf{F}) = \mathbf{I}_{m \times m}$ , and  $\text{cov}(\boldsymbol{\varepsilon}) = \boldsymbol{\Psi}$ ,

where

$$\boldsymbol{\Psi} = \begin{bmatrix} \psi_1 & 0 & \cdots & 0 \\ 0 & \psi_2 & & 0 \\ \vdots & & & \vdots \\ 0 & 0 & \cdots & \psi_p \end{bmatrix}$$

**Note:**

$$\mathbf{\Sigma} = \text{cov}(\mathbf{X}) = \mathbf{LL}' + \mathbf{\Psi}$$

Hence

$$\sigma_{ii} = \text{Var}(X_i) = \sum_{j=1}^m \lambda_{ij}^2 + \psi_i$$

and

$$\sigma_{ik} = \text{cov}(X_i, X_k) = \sum_{j=1}^m \lambda_{ij} \lambda_{kj}$$

$h_i^2 = \sum_{j=1}^m \lambda_{ij}^2$  is called the  $i^{\text{th}}$  *communality*

i.e. the component of variance of  $x_i$  that is due to the common factors  $F_1, F_2, \dots, F_m$

$\psi_i$  is called the *specific* variance

i.e. the component of variance of  $x_i$  that is **specific** only to that variable  
(sometimes  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_p$  are called the **specific factors** )

# Factor Rotation

If  $\hat{\mathbf{L}}$  is the  $p \times m$  matrix of estimated factor loadings, then let

$$\hat{\mathbf{L}}^* = \hat{\mathbf{L}}\mathbf{T}$$

where  $\mathbf{T}$  is any orthogonal matrix (that is,  $\mathbf{T}'\mathbf{T} = \mathbf{T}\mathbf{T}' = \mathbf{I}$ ).

Then  $\hat{\mathbf{L}}^*$  is a  $p \times m$  matrix of *rotated* loadings. Moreover,

$$\text{cov}(\mathbf{X}) = \hat{\mathbf{L}}\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}} = \hat{\mathbf{L}}\mathbf{T}\mathbf{T}'\hat{\mathbf{L}}' + \hat{\mathbf{\Psi}} = \hat{\mathbf{L}}^*\hat{\mathbf{L}}^{*'} + \hat{\mathbf{\Psi}}$$

That is, the estimated covariance matrix remains unchanged!

Thus, from a mathematical point of view it is the same whether  $\hat{\mathbf{L}}$  or  $\hat{\mathbf{L}}^*$  is used, because the factor model is overparametrized and has many solutions.

# Notes:

- Since the original loadings may not be easily interpretable, it is a common practice to rotate them until a “simpler structure” is achieved.
- Ideally, we want to see a pattern of loadings such that each variable loads heavily on a single factor and has small loadings on the remaining factors.
- It is not always possible to obtain such simple structure ☹️
- There are graphical and analytical methods for choosing the optimal rotation.



## **Example 1:**    *Olympic decathlon Scores*

Data was collected for  $n = 280$  starts from 1960 to 2004 for the ten decathlon events (*100-m run, Long Jump, Shot Put, High Jump, 400-m run, 110-m hurdles, Discus, Pole Vault, Javelin, 1500-m run*).

## Correlated Factors

We now change the assumption that the factors are *uncorrelated* with the following:

$$\text{cov}(\mathbf{F}) = \mathbf{\Phi} = \begin{bmatrix} \phi_{11} & \phi_{12} & \dots & \phi_{1m} \\ \phi_{21} & \phi_{22} & \dots & \phi_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{m1} & \phi_{m2} & \dots & \phi_{mm} \end{bmatrix}$$

Then the covariance of the data becomes

$$\mathbf{\Sigma} = \text{cov}(\mathbf{X}) = \mathbf{L}\mathbf{\Phi}\mathbf{L}' + \mathbf{\Psi}$$

# Structural Equation models:

## Path Diagrams

The above equation is just one way to specify the model we want to fit to the data. Alternatively, we may use *path diagrams* to specify the model graphically.

In a path diagram, observed variables are denoted with rectangles and latent (unobservable) variables (factors) are denoted with ovals.

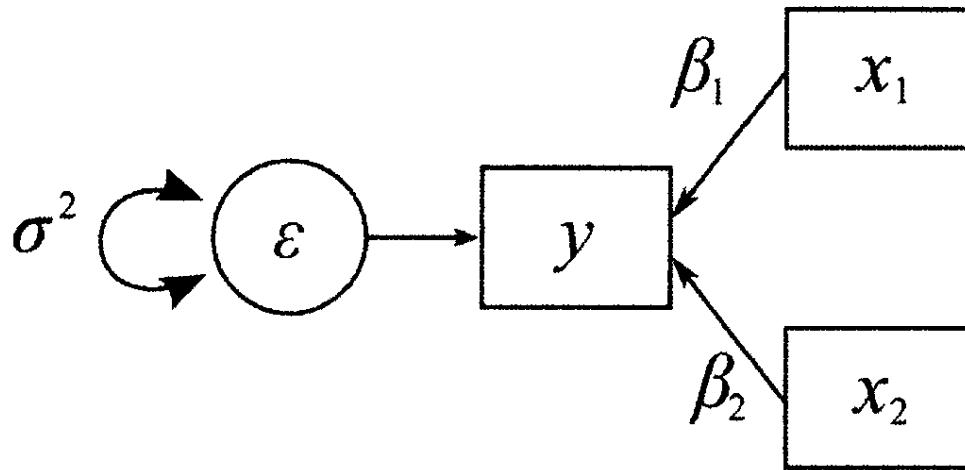
The ovals and rectangles are connected with arrows. Note that a one-headed arrow denotes a structural or unidirectional relationship, while a two-headed arrow denotes a covariance among variables. A two-headed arrow pointing from one variable and back to itself is used to specify the variance of the variable, which is usually an error or "disturbance" variable.

## Example:

Consider the two-predictor regression model

$$y = \boldsymbol{\beta}_1 x_1 + \boldsymbol{\beta}_2 x_2 + \boldsymbol{\varepsilon}, \text{Var}(\boldsymbol{\varepsilon}) = \boldsymbol{\sigma}^2$$

It can also be represented with a path diagram as follows:



## **Example 2: 14.2.2 from Rencher**

Table 14.1, which was obtained from a university business statistics class. As shown in Table 14.1, each of the 94 students in the class received scores for laboratory assignments (Lab), homework assignments (HW), pop quizzes (PopQuiz), midterm exam #1 (Exam1), midterm exam #2 (Exam2), and the final exam (FinalExam).

The instructor might hypothesize that the 6-dimensional measure of performance in the class is being driven by an underlying 2-dimensional factor process, with the first factor associated with daily effort and the second factor associated with knowledge mastery.

A possible model is:

$$\text{Lab} = \boldsymbol{\mu}_1 + \boldsymbol{\lambda}_{11}f_1 + \boldsymbol{\lambda}_{12}f_2 + \boldsymbol{\varepsilon}_1$$

$$\text{HW} = \boldsymbol{\mu}_2 + f_1 + \boldsymbol{\varepsilon}_2$$

$$\text{PopQuiz} = \boldsymbol{\mu}_3 + \boldsymbol{\lambda}_{31}f_1 + \boldsymbol{\lambda}_{32}f_2 + \boldsymbol{\varepsilon}_3$$

$$\text{Exam1} = \boldsymbol{\mu}_4 + \boldsymbol{\lambda}_{41}f_1 + \boldsymbol{\lambda}_{42}f_2 + \boldsymbol{\varepsilon}_4$$

$$\text{Exam2} = \boldsymbol{\mu}_5 + \boldsymbol{\lambda}_{51}f_1 + \boldsymbol{\lambda}_{52}f_2 + \boldsymbol{\varepsilon}_5$$

$$\text{Final Exam} = \boldsymbol{\mu}_6 + f_2 + \boldsymbol{\varepsilon}_6$$

**Note:** Some of the loadings are assumed to be 1 because of model identifiability issues.

## Example (continued):

We will further simplify the model to test the theory that variables mostly load on one of the factors:

$$\text{Lab} = \boldsymbol{\mu}_1 + \boldsymbol{\lambda}_{11}f_1 + \boldsymbol{\varepsilon}_1$$

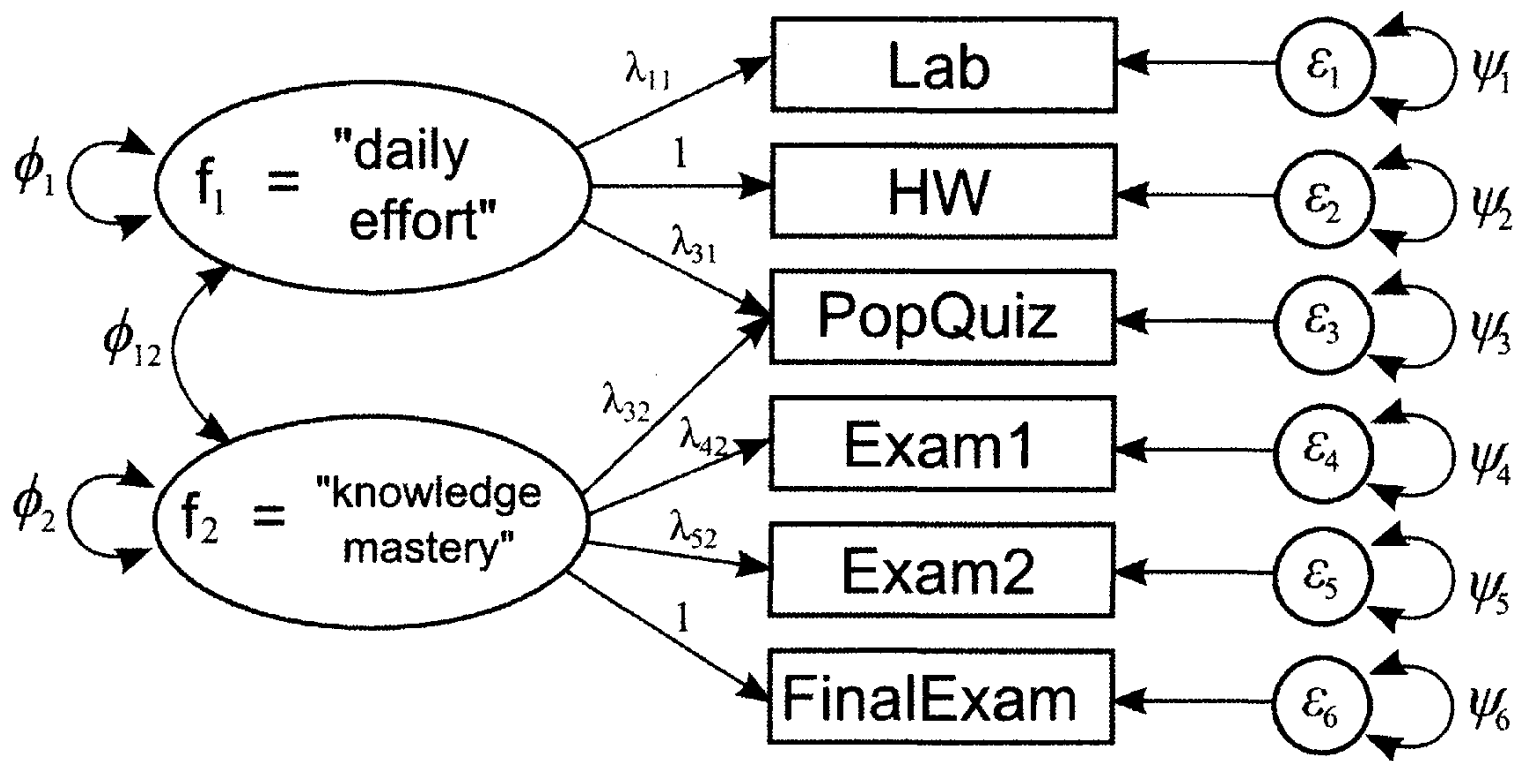
$$\text{HW} = \boldsymbol{\mu}_2 + f_1 + \boldsymbol{\varepsilon}_2$$

$$\text{PopQuiz} = \boldsymbol{\mu}_3 + \boldsymbol{\lambda}_{31}f_1 + \boldsymbol{\lambda}_{32}f_2 + \boldsymbol{\varepsilon}_3$$

$$\text{Exam1} = \boldsymbol{\mu}_4 + \boldsymbol{\lambda}_{42}f_2 + \boldsymbol{\varepsilon}_4$$

$$\text{Exam2} = \boldsymbol{\mu}_5 + \boldsymbol{\lambda}_{52}f_2 + \boldsymbol{\varepsilon}_5$$

$$\text{Final Exam} = \boldsymbol{\mu}_6 + f_2 + \boldsymbol{\varepsilon}_6$$





# Introduction to multivariate latent

- The regression model from last time is somewhat limiting because it only describes the conditional distribution of one variable given the others.
- In general, we may be interested in the relationships among all the variables in a dataset.
- If the variables were approximately jointly normally distributed, then we could describe the relationships among the variables with the sample covariance matrix or a multivariate normal model.
- However, such a model is inappropriate for nonnumeric ordinal variables.
- To accommodate ordinal variables, we can extend the ordered probit model above to a latent, multivariate normal model

# The Gaussian copula model

Let the data be  $p$ -dimensional vectors satisfying:

$$Y_{ij} = g_j(Z_{ij}), i = 1, \dots, n$$
$$\mathbf{Z}_1, \dots, \mathbf{Z}_n \sim N_p(0, \mathbf{\Psi})$$

We assume  $g_1, \dots, g_p$  are non-decreasing functions and  $\mathbf{\Psi}$  is a correlation matrix with diagonal elements equal to 1. Under this model the cdf of  $j^{\text{th}}$  component of  $\mathbf{Y}$  is  $F_j(y) = \Phi\left(g_j^{-1}(y)\right)$ .

This is known as the *multivariate normal copula model*. The term “copula” refers to the method of “coupling” a model for multivariate dependence to a model for the marginal distributions of the data.

# Copula estimation

The unknown parameters in the copula model are the matrix  $\Psi$  and the non-decreasing functions  $g_1, \dots, g_p$ . Bayesian inference for all of these parameters would require that we specify a prior for  $\Psi$  as well as  $p$  prior distributions over the complicated space of arbitrary non-decreasing functions. There is a rank likelihood method which avoids this. Since each  $g_j$  is non-decreasing, observing the  $n \times p$  data matrix  $Y$  tells us that the matrix of latent variables  $Z$  must lie in the set

$$R(Y) = \{Z: z_{i_1 j} < z_{i_2 j} \text{ if } y_{i_1 j} < y_{i_2 j}\}$$

As a function of  $\Psi$ ,  $P(Z \in R(Y) | \Psi)$  is called the rank likelihood for the multivariate normal copula model. Computing the likelihood for a given value of  $\Psi$  is very difficult, but we can make an MCMC approximation to  $f(\Psi, Z | Z \in R(Y))$  using Gibbs sampling, provided we use a prior for  $\Psi$  based on the inverse-Wishart distribution.

# A parameter-expanded prior distribution for $\Psi$

Unfortunately, there is no simple conjugate class of prior distributions for our correlation matrix  $\Psi$ . Let's consider an alternative model:

$$Y_{ij} = g_j(Z_{ij}), i = 1, \dots, n$$
$$\mathbf{Z}_1, \dots, \mathbf{Z}_n \sim N_p(0, \Sigma)$$

where  $\Sigma$  is a covariance matrix. In this case a natural prior distribution for  $\Sigma$  would be an inverse-Wishart distribution, which would give an inverse-Wishart full conditional distribution and thus make posterior inference available via Gibbs sampling. Careful inspection of the rank likelihood indicates that it does not provide us with a complete estimate of  $\Sigma$ . For this reason, we say that the diagonal entries of  $\Sigma$  are non-identifiable in this model. The value of  $\Psi$  is identifiable from the rank likelihood, and so one estimation approach for the Gaussian copula model is to reparametrize the model in terms of a non-identifiable covariance matrix  $\Sigma$ , but focus our posterior inference on the identifiable correlation matrix  $\Psi = h(\Sigma)$ .