

STAT 3690 Lecture Note

Part VIII: Factor analysis

Zhiyang Zhou (zhiyang.zhou@umanitoba.ca, zhiyanggeezhou.github.io)

2023/Mar/27 01:38:28

Factor analysis

Latent variable model

- latent/unobservable variables give rise to observed data through a specific model, i.e., a regression model with unobservable covariates
- Factor analysis model is a special kind of latent variable model

Population version

- Model

$$\mathbf{Y} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \mathbf{E}$$

- $\mathbf{Y} = [Y_1, \dots, Y_p]^\top \sim (\boldsymbol{\mu}, \boldsymbol{\Sigma})$: random & observable p -vector
- $\mathbf{L} = [\ell_{ij}]_{p \times q}$: fixed & unknown, a matrix of factor loadings
 - * ℓ_{ij} : the contribution of j th factor to Y_i
- $\mathbf{F} \sim (\mathbf{0}, \mathbf{I}_q)$: random & unobservable, q -vector of latent/common factors
- $\mathbf{E} \sim (\mathbf{0}, \boldsymbol{\Psi})$: random & unobservable, p -vector of error/specific factors, with $\boldsymbol{\Psi} = \text{diag}(\psi_1, \dots, \psi_p)$ and $\text{cov}(\mathbf{F}, \mathbf{E}) = \mathbf{0}$
- Covariance structure
 - $\text{var}(\mathbf{Y}) = \boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^\top + \boldsymbol{\Psi}$
 - * I.e., $\text{var}(Y_i) = \sum_{j=1}^q \ell_{ij}^2 + \psi_i$
 - $\text{cov}(\mathbf{Y}, \mathbf{F}) = \mathbf{L}$
 - $\sum_{i=1}^p \ell_{ij}^2$: the variability contributed by the j th latent factor

Sample version

- Model

$$\mathbf{Y}_i - \boldsymbol{\mu} = \mathbf{L}\mathbf{F}_i + \mathbf{E}_i, \quad i = 1, \dots, n$$

- $\mathbf{Y}_1, \dots, \mathbf{Y}_n \stackrel{\text{iid}}{\sim} \mathbf{Y}$
- $\mathbf{F}_1, \dots, \mathbf{F}_n \stackrel{\text{iid}}{\sim} \mathbf{F}$
- $\mathbf{E}_1, \dots, \mathbf{E}_n \stackrel{\text{iid}}{\sim} \mathbf{E}$

Estimating \mathbf{L} and $\boldsymbol{\Psi}$

- Selection of q , i.e., the number of latent factors, with one of the following rules
 - PCA stopping rule

- Taking q such that $\sum_{j=1}^q \sum_{i=1}^p \ell_{ij}^2 / \text{tr}(\mathbf{S})$ is over a preset percentage, where $\mathbf{S} = (n-1)^{-1} \sum_{i=1}^n (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})^\top$
 - * $\sum_{i=1}^q \ell_{ij}^2 / \text{tr}(\mathbf{S})$: the proportion of variation explained by the j th latent factor
- Taking q as the number of positive eigenvalues of \mathbf{S}
- Taking q as the number of eigenvalues of \mathbf{S} that are above average
- Taking q as the number of eigenvalues of correlation matrix greater than one
- According to domain-knowledge expertise
- PC method
 1. Determine q
 2. Pick up the q largest eigenvalues $\lambda_1, \dots, \lambda_q$ of \mathbf{S} and corresponding eigenvectors w_1, \dots, w_q
 3. $\hat{\mathbf{L}} = [\sqrt{\lambda_1}w_1, \dots, \sqrt{\lambda_q}w_q]_{p \times q}$ and $\hat{\Psi} = \text{diag}(\mathbf{S} - \hat{\mathbf{L}}\hat{\mathbf{L}}^\top)$

-
- Exercise 8.1: `psych::bfi` involves 2800 subjects, covering their 25 personality assessments, gender, education and age.

```
install.packages(c('psych'))
library(psych)
library(tidyverse)
options(digits = 4)
head(psych::bfi)
data = bfi %>%
  select(-gender, -education, -age) %>% # Remove gender, education and age
  filter(complete.cases(.)) # keep complete data only
S = cov(data)
fa_pc = prcomp(data) # decompose the covariance matrix

# PCA stopping rule
(q = which(cumsum(fa_pc$sdev^2)/sum(fa_pc$sdev^2)>.9)[1])
# the overall proportion of variation explained by latent factors
(q = which(
  cumsum(sort(colSums((fa_pc$rotation %*% diag(fa_pc$sdev))^2), decreasing = T))/sum(diag(S))>.9
)[1])
# the number of eigenvalues above the average
(q = sum(eigen(S, only.values = T)$values > mean(eigen(S, only.values = T)$values)))
# the number of eigenvalues greater than one for the correlation matrix
(q = sum(eigen(cor(data))$values > 1))

L_pc = fa_pc$rotation[,1:q] %*% diag(fa_pc$sdev[1:q])
Psi_pc = diag(diag(S - tcrossprod(L_pc)))

S_pc = tcrossprod(L_pc) + Psi_pc
lattice::levelplot(S - S_pc, scales=list(x=list(rot=90))) # fitting error
lattice::levelplot((S - S_pc)/S, scales=list(x=list(rot=90))) # difference in percentage
```

-
- ML method
 - Assuming
 - * $\mathbf{F} \sim \text{MVN}_q(\mathbf{0}, \mathbf{I})$
 - * $\mathbf{E} \sim \text{MVN}_p(\mathbf{0}, \Psi)$
 - * Diagonal $\mathbf{L}^\top \Psi^{-1} \mathbf{L}$
 - Resorting to R functions `factanal` or `psych::fa`
-

```

install.packages(c('psych'))
library(psych)
library(tidyverse)
options(digits = 4)
head(psych::bfi)
data = bfi %>%
  select(-gender, -education, -age) %>% # Remove gender, education and age
  filter(complete.cases(.)) # keep complete data only
S = cov(data)

# the number of eigenvalues greater than one for the correlation matrix
(q = sum(eigen(cor(data))$values > 1))

# apply functions factanal OR psych::fa
fa_ml_1 <- factanal(covmat = S, factors = q, rotation = 'none')
fa_ml_2 <- psych::fa(r = S, covar = T, nfactors = q, rotate = "none", fm = "ml")
head(fa_ml_1$loadings-fa_ml_2$loadings)

L_ml <- fa_ml_1$loadings
Psi_ml <- diag(fa_ml_1$uniquenesses)
S_ml = tcrossprod(L_ml) + Psi_ml
lattice::levelplot(S - S_ml,
  scales=list(x=list(rot=90)), xlab = "", ylab = "")
lattice::levelplot((S - S_ml)/S,
  scales=list(x=list(rot=90)), xlab = "", ylab = "")

```

-
- Comments on the estimation of \mathbf{L} and $\mathbf{\Psi}$
 - More methods other than ML and PC
 - Different statistical softwares may apply different methods
 - * Have to look into help manuals to figure out what is going on for different softwares/packages
 - Compare the outputs of multiple estimation methods
 - * For a good fit, similar answers would be reached regardless of the method
-

Factor rotation

- \mathbf{L} is not uniquely defined: if $\mathbf{Y} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \mathbf{E}$, then $\mathbf{Y} - \boldsymbol{\mu} = \tilde{\mathbf{L}}\tilde{\mathbf{F}} + \mathbf{E}$, where
 - $\tilde{\mathbf{L}} = \mathbf{L}\mathbf{P}$ and $\tilde{\mathbf{F}} = \mathbf{P}^\top \mathbf{F}$ with \mathbf{P} a $q \times q$ orthogonal matrix ($\mathbf{P}^{-1} = \mathbf{P}^\top$)
- A blessing to improve interpretation: pick up a \mathbf{P} such that $\tilde{\mathbf{F}}$ is more interpretable; to ease the interpretation, we want:
 - Each entry of \mathbf{Y} to have large loadings for merely one latent factor and negligible loadings for remaining ones
- varimax: find \mathbf{P} to maximize the sum of variance of squared (scaled) loadings over all the latent factors

$$\sum_{j=1}^q \left\{ \frac{1}{p} \sum_{i=1}^p \tilde{\ell}_{ij}^{*4} - \left(\frac{1}{p} \sum_{i=1}^p \tilde{\ell}_{ij}^{*2} \right)^2 \right\}$$

$$- \tilde{\ell}_{ij}^* = \tilde{\ell}_{ij} / \sqrt{\sum_{j=1}^q \tilde{\ell}_{ij}^2} \text{ with } \tilde{\ell}_{ij} \text{ the } (i, j)\text{-th entry of } \tilde{\mathbf{L}} = \mathbf{L}\mathbf{P}$$

- Comments on factor rotation
 - Especially useful with loadings obtained through ML
 - Sometimes used even for loadings in PCA

Factor scores

- Weighted least square (WLS) method
 - Given $\bar{\mathbf{Y}}$, $\hat{\mathbf{L}}$, and $\hat{\mathbf{\Psi}}$, then, for the i th observation \mathbf{Y}_i ,

$$\hat{\mathbf{F}}_i = (\hat{\mathbf{L}}^\top \hat{\mathbf{\Psi}}^{-1} \hat{\mathbf{L}})^{-1} \hat{\mathbf{L}}^\top \hat{\mathbf{\Psi}}^{-1} (\mathbf{Y}_i - \bar{\mathbf{Y}})$$

- * I.e., the minimizer of $(\mathbf{Y}_i - \bar{\mathbf{Y}} - \hat{\mathbf{L}}\mathbf{F})^\top \hat{\mathbf{\Psi}}^{-1} (\mathbf{Y}_i - \bar{\mathbf{Y}} - \hat{\mathbf{L}}\mathbf{F})$ with respect to \mathbf{F}
-
-

- Regression method
 - Assuming $\mathbf{F} \sim \text{MVN}_q(\mathbf{0}, \mathbf{I})$ and $\mathbf{E} \sim \text{MVN}_p(\mathbf{0}, \mathbf{\Psi})$,

$$\begin{bmatrix} \mathbf{Y} - \boldsymbol{\mu} \\ \mathbf{F} \end{bmatrix} \sim \text{MVN}_{p+q} \left(\mathbf{0}, \begin{bmatrix} \mathbf{L}\mathbf{L}^\top + \mathbf{\Psi} & \mathbf{L} \\ \mathbf{L}^\top & \mathbf{I} \end{bmatrix} \right)$$

and hence

$$\mathbf{F} | \mathbf{Y} \sim \text{MVN}_p(\mathbf{L}^\top (\mathbf{L}\mathbf{L}^\top + \mathbf{\Psi})^{-1} (\mathbf{Y} - \boldsymbol{\mu}), \mathbf{I} - \mathbf{L}^\top (\mathbf{L}\mathbf{L}^\top + \mathbf{\Psi})^{-1} \mathbf{L})$$

- Given $\bar{\mathbf{Y}}$, $\hat{\mathbf{L}}$, and $\hat{\mathbf{\Psi}}$, estimate \mathbf{F}_i by

$$\hat{\mathbf{F}}_i = \hat{\mathbf{L}}^\top (\hat{\mathbf{L}}\hat{\mathbf{L}}^\top + \hat{\mathbf{\Psi}})^{-1} (\mathbf{Y}_i - \bar{\mathbf{Y}})$$

OR

$$\hat{\mathbf{F}}_i = \hat{\mathbf{L}}^\top \mathbf{S}^{-1} (\mathbf{Y}_i - \bar{\mathbf{Y}})$$

- Comments on factor scores
 - More methods available
 - No uniformly superior way

Summary on factor analysis

- What we discussed is “exploratory” factor analysis
 - “Confirmatory” factor analysis would make stronger assumptions about the nature of the latent factors and perform statistical inference.
 - There are choices to make at every stage of factor analysis: estimation method, number of factors, factor rotation, and score estimation.
 - * Too flexible to be tracked
 - * Close to an “art”
- General strategy for factor analysis
 1. Perform a PC factor analysis
 - It may help you identify potential outliers
 2. Perform an ML factor analysis.
 - Try a varimax rotation to see if it makes sense
 3. Compare the solutions of both methods to see if they generally agree.
 4. Repeat for different number of common factors q and check if adding more factors may improve the interpretation
 5. For large datasets, you can split your data, run the same model on both subsets, and compare the loadings to see if they generally agree

An example of factor analysis

- `state.x77` contains general information about all 50 US states
 - Population
 - Income per capita
 - Illiteracy rate
 - Life expectancy
 - Murder rate
 - High-school graduation rate
 - Average number of freezing degree days (with the temperature lower than 0 °C)
 - Total area
-