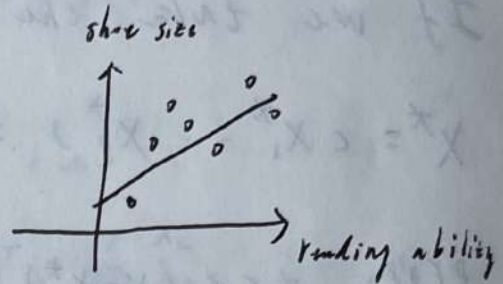# Factor Analysis

(1) Background:

i) A motivating example:

shoe size and reading ability

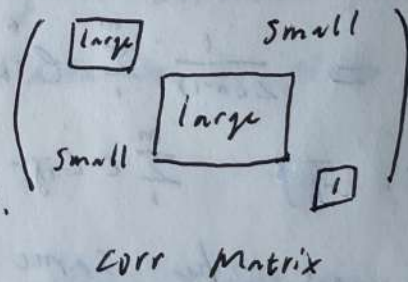exist strong correlation

$\Rightarrow$ Latent variable : age.



ii) Purpose of Factor analysis:

• Reduce high dimensional data to a few representive Variables

• Describe the relation between variables (correlation or covariance) by underlying variables.

e.g. Variables can be grouped by its correlation matrix.



corr Matrix

Rmk: PCA is different from FA. Since PCA just transform the data matrix to reduce its dimension. It doesn't need modeling. But Factor analysis needs to find the underlying communal factors.

(2) Modeling:

① Orthogonal Factor Model:

$X = (X_1 \cdots X_p)^T$ observed random vector with $\begin{cases} E(X) = \mu \\ Var(X) = \Sigma \end{cases}$.

Assumptions: $X$ is linear dependent upon a few r.v's $F_1 \cdots F_m$ called common factors and $p$ specific factors (errors) $\varepsilon_1 \cdots \varepsilon_p$. All are unobservable.

Satisfies: i) $\overline{E}(\varepsilon) = 0_{p \times 1}$. $Cov(\varepsilon) = \Psi = diag\{\psi_1 \cdots \psi_p\}$.

ii) $\overline{E}(F) = 0_{m \times 1}$  $Cov(F) = I_m$.

iii) $F$ and $\varepsilon$ are indept. $Cov(F, \varepsilon) = 0$.

Rmk: ii) isn't strong since we can do PCA to decompose it orthogonally.

Model: $X_i - \mu_i = \sum\limits_{k=1}^{m} \ell_{ik} F_k + \varepsilon_i$. $1 \leq i \leq P$.

Written in matrix: $X - \mu = L F + \varepsilon$.

Rmk: i) We want $m < P$ as possible.

ii) Likewise regression model: $Y = X\beta + \varepsilon$.

Since $L$ is invariant. it's kind of $\beta$. And note that $X$ consists different kinds of variables which differs a lot from regression model.

Def: $\ell_{ij}$ is factor loading of $i^{th}$ variable on $j^{th}$ factor.

Some results: i) $Cov(X) = L L^T + \Psi$. i.e.

$Cov(X_j, X_k) = \sum\limits_{i=1}^{m} \ell_{ji} \ell_{ki}$. $j \neq k$.

$$\sigma_{ii} = \sum_{k=1}^{m} l_{ik}^2 + \psi_i =: \text{"Commanality"} + \text{"Specific Var"}.$$

ii) $\sum_i^p \text{Var}(X_i) = \sum h_i^2 + \sum \psi_i$ . where $h_i^2 = \sum_{k=1}^m l_{ik}^2$ .

iii) $\text{Cov}(X, F) = L$ , i.e. $\text{Cov}(X_i, F_j) = l_{ij}$ .

### Interpretions:

i) $h_i^2$ can be recognized as kind of geometric distance. It can measure the influence of $F_1 \dots F_m$ on $X_i$.

ii) Set $q_j^2 = \sum_{i=1}^p l_{ij}$ . It represent the influence of $F_j$ on the whole model.

Rmk: The model holds under scale transform:

i.e. For $C = \text{diag}\{c_1, \dots, c_p\}$. $c_i \neq 0$. $X^* = CX$.

$\mu^* = C\mu$ . $\Sigma^* = C\Sigma C^T$. $\varepsilon^* = C\varepsilon$. Then:

$X^* = \mu^* + A^* F + \varepsilon^*$ still holds. $(A^* = CA)$

## ② Data Reduction:

The factor model assumes $p + \frac{p(p+1)}{2} = \frac{p(p+1)}{2}$ (w.r.t $\Sigma$)

variances. It can be reduced to $pm$ factor loadings and $p$ specific variances $\psi_i$. if $(m+1)p < \frac{p(p+1)}{2}$.

i.e. $m+1 < \frac{p+1}{2}$

Rmk: Unfortunately. Not every Covariance matrix can be written in $LL^T + \psi$ form. where $m < p$.

③ Rotation Indeterminary:

Let $T$ is orthogonal matrix. $L^* = LT$. $F^* = TF$. Then:
$X - \mu = LF + \epsilon = L^* F^* + \epsilon$. $E(F^*) = 0$, $Cov(F^*) = I_m$.
$\Sigma = LL^T + \Psi = L^* L^{*T} + \Psi$.

i.e. $L$ and $F$ aren't unique. They have same properties under rotation.

Rmk: i) Var "$\Sigma$" is uneffected : because it's some kind of distance

ii) To make the $L$, $F$ unique. We always impose some condition on them.

(3) Estimation:

· We will estimate $L$, $\Psi$ and number $m$. $\Sigma = LL^T + \Psi$ will be estimated by $S$. (i.e. $S \approx \hat{L}\hat{L}^T + \hat{\Psi}$)

① Principle Component Approach:

By spectral decomposition on $\Sigma$:
$$\Sigma = \sum_1^p \lambda_i e_i e_i^T = (\sqrt{\lambda_1} e_1 \cdots \sqrt{\lambda_p} e_p) \begin{pmatrix} \sqrt{\lambda_1} e_1^T \\ \vdots \\ \sqrt{\lambda_p} e_p^T \end{pmatrix} = LL^T$$
where we assume: $\lambda_1 \geq \cdots \geq \lambda_p \geq 0$.

If the last $p-m$ eigenvalues are small. Then neglect them. $\Sigma \approx (\sqrt{\lambda_1} e_1 \cdots \sqrt{\lambda_m} e_m) \begin{pmatrix} \sqrt{\lambda_1} e_1^T \\ \vdots \\ \sqrt{\lambda_m} e_m^T \end{pmatrix} = \breve{L}\breve{L}^T$

Specific variance may be taken $diag\{\Sigma - \breve{L}\breve{L}^T\}$.

Rmk: $\Sigma - \breve{L}\breve{L}^T$ may not be diagonal matrix.

Next. decompose $S = \sum_i^p \hat{\lambda}_i \hat{e}_i \hat{e}_i^T$. And set:

$$\hat{L} = (\sqrt{\hat{\lambda}_1} \, \hat{e}_1 \hat{e}_1^T - \cdots \sqrt{\hat{\lambda}_m} \, \hat{e}_m \hat{e}_m^T)$$

i) $\psi$ is estimate by:

$$\hat{\psi} = \begin{pmatrix} \hat{\psi}_1 & & 0 \\ & \ddots & \\ 0 & & \hat{\psi}_p \end{pmatrix}. \qquad \hat{\psi}_i = S_{ii} - \sum_j^m \hat{\ell}_{ij}^2$$

ii) Communalities $\tilde{h}_i$ is estimate by:

$$\hat{h}_i^2 \simeq \sum_{k=1}^m \hat{\ell}_{ik}^2.$$

Rmk: In this approach. It doesn't change the estimated loadings of given factors as the number of factors increase $m$ to $m+1$.

iii) Select number $m$:

Consider $S - (\hat{L}\hat{L}^T + \hat{\psi}) = \begin{pmatrix} 0 & & * \\ & \ddots & \\ * & & 0 \end{pmatrix} =: E_s$

$tr(E_s^T E_s) =$ sum of square entries of $E_s$

$$\leq tr((S - \hat{L}\hat{L}^T)^T (S - \hat{L}\hat{L}^T)) = \text{SS of } S - \hat{L}\hat{L}^T.$$

$$= \sum_{m+1}^p \hat{\lambda}_k^2 \qquad (\text{Apply the spectral decompose})$$

$\Rightarrow$ Select $m$. s.t. $\sum_{m+1}^p \hat{\lambda}_k^2$ is small enough.

(which reduce the SSE of estimate)

Rmk: Alternatively. Set $P_0$. s.t. $\dfrac{\sum \lambda_i}{\sum_i \lambda_i} \geq P_0$.

Besides, we can consider the proportion of $j^{th}$ factor $\hat{\lambda}_j / \sum_i^p S_{ii}$. to select.

① Modified PCA: Principle Factor Solution:

---

Ideal: The common factors should account for off diagnoal as well as community portion of diagnoal elements of $\Sigma$.

Algorithm: i) Guess $\hat{\Psi}$.

ii) Set $L$ = the largest $m$ eigenvectors of decomposition of $S - \hat{\Psi}$.

iii) Set $\hat{\Psi} = \text{diag}(S - LL^T)$.

Repeat ii) and iii) untill convergence.

Rmk: i) The first step of estimate on $\hat{\Psi}$ and then set $S - \hat{\Psi}$ to decompose can ease the estimate of $L$ on the diagnoal elements. So $L$ can better estimate the off-diagnoal part of $\Sigma$.

ii) Choice of $\hat{\Psi}$ can be $S^{-1}$.

iii) $S - \hat{\Psi}$ may have negative eigenvalue. and estimate of $\hat{\Psi}$ in step iii) may also have negative diagnoal. It's referred to Heywood case.

② Maximum Likelihood Method:

---

Assumption: i) F and $\varepsilon$ are jointly nomal distribution.

ii) To make $L$ well-def and uniqueness:

$$L^T \psi^{-1} L = \Delta \text{ is diagonal matrix.}$$
$$(\Leftrightarrow (\psi^{-\frac{1}{2}} L)^T (\psi^{-\frac{1}{2}} L) = \Delta )$$

Rmk: ii) put $\frac{m(m-1)}{2}$ constraints to reduce the dimension of para. space to 1.

$\left( \overset{\frac{m(m-1)}{2}}{\underset{\nabla}{\oslash}} @ \right)$

$$\Rightarrow L(M, \Sigma | X) = (2\pi)^{-\frac{np}{2}} |\Sigma|^{-\frac{n}{2}} \exp(-\frac{1}{2} tr( \Sigma^{-1} \sum_{i}^{n} (X_i - M)(X_i - M)^T )$$

To maximize it $\Leftrightarrow$ maximize $\ln|\Sigma| - \ln|S_n| + tr( \Sigma^{-1} S_n) - p$

where $\Sigma = LL^T + \psi$. $S_n = \frac{(n-1)S}{n}$. subjects to: $L^T \psi^{-1} L = \Delta$

MLE of $L$. $\psi$ can be obtained by numerical Computation.

Rmk. By invariant property of MLE. The MLE of communalities are $\hat{h_i}^2 = \sum_{k=1}^{m} \hat{l}_{ik}^2$ of $\hat{L}$.

For Standardization:

Set $Z = V^{-\frac{1}{2}} (X - M)$. So $\ell = V^{-\frac{1}{2}} \Sigma V^{-\frac{1}{2}}$. Set $L_z = V^{-\frac{1}{2}} L$.

$\psi_z = V^{-\frac{1}{2}} \psi V^{-\frac{1}{2}}$. Then $\ell = L_z L_z^T + \psi_z$. $(V = \begin{pmatrix} \sigma_n & 0 \\ 0 & \sigma_{pp} \end{pmatrix})$

By invariant property: $\hat{\ell} = (\hat{V}^{-\frac{1}{2}} \hat{L})(\hat{V}^{-\frac{1}{2}} \hat{L})^T + \hat{V}^{-\frac{1}{2}} \hat{\psi} \hat{V}^{-\frac{1}{2}}$.

where $\hat{V}^{-\frac{1}{2}}$. $\hat{L}$ are MLE of $V^{-\frac{1}{2}}$. $L$.

Rmk: i) For PCA. There're usually no relation between PC of "$\Sigma$" and "$R$". But in this method. They're essentially equi.

ii) The result will be very different when increase the number from $m$ to $m+1$.

That's because we impose different assumption comparing to PCA method.

iii) Heywood case may also happen ($\hat{\psi}_i < 0$).

iv) MLE method produce $F_i$ won't be orthogonal:

$$(\underline{\psi}^{-\frac{1}{2}} L)^T (\underline{\psi}^{-\frac{1}{2}} L) = \Delta$$

④ <u>Large Sample Test for $m$:</u>

<u>Assumption:</u> i) $F$ and $\varepsilon$ are jointly normal distribution.

ii) $L^T \psi^{-1} L = \Delta$ diagnoal matrix.

Test: $H_0 : \Sigma_{p\times p} = L L^T + \psi_{p\times p}$. $L \in M^{p\times m}$ v.s. $H_1 : \Sigma$ any other.

i) $\sup L(\Sigma, M) = L(S_n, \bar{X})$

ii) $\sup L(\Sigma, M) = L(\hat{L}\hat{L}^T + \hat{\psi}, \bar{X})$. $\hat{L}, \hat{\psi}$ are MLE of $H_0$.

$L$ and $\psi$. $L(\hat{L}\hat{L}^T + \hat{\psi}, \bar{X}) \propto |\hat{L}\hat{L}^T + \hat{\psi}|^{-\frac{n}{2}} \exp(-\frac{n}{2} tr((\hat{L}\hat{L}^T + \hat{\psi})^{-1} S_n))$

rmk: $tr((\hat{L}\hat{L}^T + \hat{\psi})^{-1} S_n) = p$.

$\Rightarrow -2 \ln \Delta = n \ln \left( \frac{|\hat{L}\hat{L}^T + \hat{\psi}|}{|S_n|} \right) \sim \chi^2(df)$.

$df = \nu - \nu_0 = \frac{1}{2} p(p+1) - [p(m+1) - \frac{1}{2} m(m-1)]$

$\qquad\qquad = \#$ para. of $\Sigma - [\#$ para. of $L, \psi - \#$ constraint$]$

$\qquad\qquad = \frac{1}{2}[(p-m)^2 - p - m]$

<u>rmk:</u> For $df \geq 0 \Rightarrow m < \frac{1}{2}(2p+1 - \sqrt{8p+1})$

Bartlett show the converge of approxi. can

be improved if replace $n$ by $n-1-(2p+4m+5)/6$

$$\Rightarrow R = \{ (n-1 - \frac{2p+4m+5}{4}) \ln \frac{|\hat{L}\hat{L}^T + \hat{\Psi}|}{|S_n|} > \chi^2 df \ (4) \}.$$

## (4) Factor Rotation:

Since the original loading may not be

readily interpretable. e.g. some $l_{ij}$ are positive
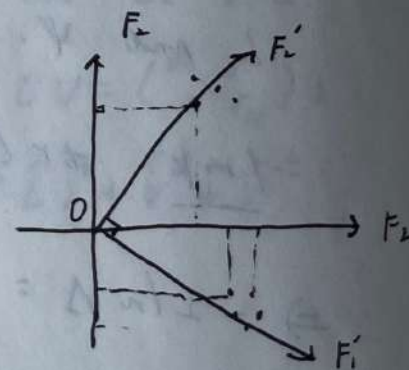
and some are negative which are large.

To achieve a simple structure. We need

rotate $L$ and $F$.

Rmk: Dimension of rotation for factor loading

is $\frac{m(m-1)}{2}$. For $T = (T_1 \cdots T_m)$. Constraint:

$T_1 \perp T_2 ; \ T_1 \perp T_3 . \ T_2 \perp T_3 , \ \cdots \ 1+2+ \cdots +(m-1).$

### ① Simple Structure:

Ideally. We like to see a pattern

that each variable loads highly on

a single factor and loads small

on the remaining factors.

e.g. Rotate $F_1, \circ F_2$ to $F_1', \circ F_2'$. A partition

into mutually exclusive groups would be

desirable.

① Varimax Criteria:

· Suppose $\hat{L}^* = (\hat{l}_{ij}^*)$ is $\hat{L}$ after rotation

Def: $\tilde{l}_{ij}^* = \hat{l}_{ij}^* / \hat{h}_i$. (i.e. scaling, consider the weight)

$\tilde{l}_{ij}^2 = \tilde{l}_{ij}^{*2}$. $\bar{\lambda}_j = \sum_{i=1}^{p} \tilde{l}_{ij}^2 / p$.

$\Rightarrow V_j = \sum_{i=1}^{p} (\tilde{l}_{ij}^2 - \bar{\lambda}_j)^2 / p = \frac{1}{p} \left[ \sum_{i=1}^{p} \frac{\hat{l}_{ij}^{*4}}{\hat{h}_i^4} - \left( \sum_{i=1}^{p} \frac{\hat{l}_{ij}^{*2}}{\hat{h}_i^2} \right)^2 / p \right]$

$V = \sum_{j=1}^{m} V_j = \frac{1}{p} \sum_{j=1}^{m} \left[ \sum_{i=1}^{p} \tilde{l}_{ij}^{*4} - \left( \sum_{i=1}^{p} \tilde{l}_{ij}^{*2} \right)^2 / p \right]$

Select $T$ st. make $V$ as large as possible. i.e. make the data loading of $F_j$ disperse enough.

Interpretation: $V \propto \sum_{j=1}^{m} \left( \begin{array}{l} \text{Var of square of scaled} \\ \text{loading for } j^{th} \text{ factor} \end{array} \right)$
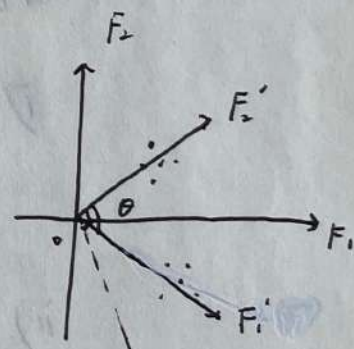
Rmk: In some case, even orthogonal rotate don't provide an easy interpretation. It's possible to use the oblique rotations at expense of orthogonality of factors. i.e.
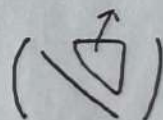
$Cov(F^*) \neq I_m$



e.g. $m=2$, $A = \begin{pmatrix} a_{11} & a_{12} \\ \vdots & \vdots \\ a_{p1} & a_{p2} \end{pmatrix}$. Set $T = \begin{pmatrix} \cos\gamma & -\sin\gamma \\ \sin\gamma & \cos\gamma \end{pmatrix}$

$B = AT = \begin{pmatrix} b_{11} & b_{12} \\ \vdots & \vdots \\ b_{p1} & b_{p2} \end{pmatrix}$ $V_* = \frac{1}{p^2} \left[ p \sum_{i=1}^{p} b_{i*}^4 / h_i^4 - \left( \sum_i \frac{b_{i*}^2}{h_i^2} \right)^2 \right]$

$\frac{\partial V}{\partial \gamma} = \frac{\partial (V_1 + V_2)}{\partial \gamma} = 0 \Rightarrow$ obtain $\gamma$

For $m > 2$, consider $(F_i, F_j)$, $i \neq j$, $\binom{m}{2}$ times

(2)

(★) Factor Scores:

We will predict/estimate the unobservable

random factors $F_1, F_2 \cdots F_m$.

① Weighted Leasted Square:

$X = \mu + LF + \varepsilon$. $\varepsilon \sim N(0, \Psi)$ Assume. $L, \Psi, \mu$ are known.

minimize: $\varepsilon^T \Psi^{-1} \varepsilon = (X - \mu - LF)^T \Psi^{-1} (X - \mu - LF)$

$\Rightarrow$ Weighted LSE is $\hat{F} = (L^T \Psi^{-1} L)^{-1} L^T \Psi^{-1} (X - \mu)$

Rmk: Practically. $L, \Psi, \mu$ are unknown. We

take estimate $\hat{L}, \hat{\Psi}, \hat{\mu} = \bar{X}$ to replace

it in WLSE. $\Rightarrow \hat{F}_j = (\hat{L}^T \hat{\Psi}^{-1} \hat{L})^{-1} \hat{L}^T \hat{\Psi}^{-1} (X_j - \bar{X})$

i) MLE. method:

$\hat{L}^T \hat{\Psi}^{-1} \hat{L} = \hat{\Delta} \Rightarrow \hat{F}_j = \hat{\Delta}^{-1} \hat{L}^T \hat{\Psi}^{-1} (X_j - \bar{X})$

ii) Correlation matrix factors:

$\hat{F}_j = (\hat{L}_z^T \hat{\Psi}_z^{-1} \hat{L}_z)^{-1} \hat{L}_z^T \hat{\Psi}_z^{-1} Z_j$. $Z_j = V^{-\frac{1}{2}} (X_j - \bar{X})$

② Regression Method:

i) Thompson Factors:

Express F in X: $F_i = \sum\limits_{k=1}^{\prime} \beta_{ik} (X_k - \mu) = \sum\limits' \beta_{ik} \bar{X}_k$.

Develop model: $F_i = \sum\limits_{k=1}^{\prime} b_{ik} \bar{X}_k$. $1 \leq i \leq m$.

$\Rightarrow \ell_{ij} = Cov(F_j, X_i) = \sum\limits_{k=1}^{\prime} b_{jk} \sigma_{ik}$.

$\therefore B = L^T \Sigma^{-1}$. We obtain: $\hat{F} = L^T \Sigma^{-1}(X-M)$.

$$\Rightarrow \hat{F} = \hat{L}^T(\hat{L}\hat{L}^T + \hat{\Psi})^{-1}(X-\bar{X}).$$

ii) Bayesian Method:

Assume $F, \epsilon$ are jointly normal distribution.

$$\begin{pmatrix} X-M \\ F \end{pmatrix} \sim N_{m+p}(0, \Sigma^*), \quad \Sigma^* = \begin{pmatrix} \Sigma & L \\ L^T & I \end{pmatrix}, \quad \Sigma = LL^T + \Psi.$$

$$\Rightarrow \begin{cases} E(F|X) = \bar{E}(F) + L^T \Sigma^{-1}(X-M) = L^T \Sigma^{-1}(X-M) \\ Cov(F|X) = I - L^T \Sigma^{-1} L. \end{cases}$$

$\therefore \hat{F} = \hat{L}^T(\hat{L}^T\hat{L} + \hat{\Psi})^{-1}(X-\bar{X})$. identical with i).

③ Comparison:

Note that: $\hat{L}^T(\hat{L}\hat{L}^T + \hat{\Psi})^{-1} = (I + \hat{L}^T\hat{\Psi}^{-1}\hat{L})^{-1}\hat{L}^T\hat{\Psi}^{-1}$.

(By calculate $\begin{pmatrix} \hat{\Psi} & \hat{L} \\ \hat{L}^T & I \end{pmatrix}^{-1}$)

$$\Rightarrow \hat{F}^{LS} = (\hat{L}^T\hat{\Psi}^{-1}\hat{L})^{-1}(I + \hat{L}^T\hat{\Psi}^{-1}\hat{L})\hat{F}^R$$

$$= (I + (\hat{L}^T\hat{\Psi}^{-1}\hat{L})^{-1})\hat{F}^R$$

Rmk: In MLE method, if $\hat{A} = (\hat{L}^T\hat{\Psi}\hat{L}) \approx 0$.

then $\hat{F}^{LS} \approx \hat{F}^R$.

i) $E(\hat{F}^{LS}|F) = F$. unbiased.

$E(\hat{F}^R|F) = (I + L^T\Psi^{-1}L)^{-1}L^T\Psi^{-1}LF$. biased.

ii) $E((\hat{F}^{LS}-F)(\hat{F}^{LS}-F)^T) = (L^T\Psi L)^{-1}$.

$E((\hat{F}^R-F)(\hat{F}^R-F)^T) = (I + L^T\Psi^{-1}L)^{-1}$.

(6) Strategies:

   i) Perform PCA

   ii) Perform MLE

   iii) Compare Solution of i). ii).

   iv) Repeat i), ii), iii) for other number of common factors $m$.

   v) For large data sets. split them into half and perform FA.

Rmk: i) To reduce the effect of incorrect determination of number $m$. $S$ is often used to estimate $\hat{\Sigma}$ rather than $\hat{L}\hat{L}^T + \hat{\Psi}$.

   ii) After rotation $\hat{L}^* = T\hat{L}$. Then $\hat{F}^* = T^T\hat{F}$. should be used.

   iii) After rotation, difference between estimate of MLE and PCA may be little. Since it breaks the constraints.