13/04/15.

MLA

# FACTOR ANALYSIS

- Factor Analysis is a mathematical approach for attempting to explain the correlation between a large set of variables in terms of a small number of underlying factors

- Primary assumption of FA is that we cannot observe these factors directly → they are latent. It is a dimension reduction technique, like PCA except more elaborate

- Example: correlation matrix for performance of students in Irish $(x_1)$, english$(x_2)$, math$(x_3)$

$$\begin{pmatrix} 1 & 0.8 & 0.6 \\ 0.8 & 1 & 0.85 \\ 0.6 & 0.85 & 1 \end{pmatrix}$$

The dimensionality of matrix can be reduced from $m=3$ to $m=1$ by expressing the three variables as

$$x_1 = \lambda_1 f + \varepsilon_1 \,, \quad x_2 = \lambda_2 f + \varepsilon_2 \quad x_3 = \lambda_3 f + \varepsilon_3$$

- The f's in these equations are an underlying common factor which can often be given an interpretation (in example it could be general ability)

- The $\lambda_i$'s are called factor loadings, and the $\varepsilon_i$'s are errors / specific factors

- $\varepsilon_i$'s will be of small variance if $x_i$ is closely related to f.

- More generally we have observable random vector $X = (X_1 \ldots X_n)$ with mean $\mu$ and covariance matrix $\Sigma$.

- The factor model states that X is linearly dependent upon a few unobservable random variables $f_1, \ldots f_p$ called common factors, and m additional sources of variation $\varepsilon_1, \ldots \varepsilon_m$ called specific factors

We have 
$$X_1 - \mu_1 = \lambda_{11} f_1 + \lambda_{12} f_2 + \ldots \lambda_{1p} f_p + \varepsilon_1$$
$$\vdots$$
$$X_m - \mu_m = \lambda_{m1} f_1 + \lambda_{m2} f_2 + \ldots + \lambda_{mp} f_p + \varepsilon_m$$
$$\Rightarrow X - \mu = \Lambda f + \varepsilon$$

- The this $\lambda_{ij}$ value is called the factor loading of the $i^{th}$ variable on the $j^{th}$ factor

- $\Lambda$ is the matrix of factor loading
- Note that the $i^{th}$ specific factor $\varepsilon_i$ is associated only with the response $X_i$
- Note that $f_1, f_2, \ldots f_p$ $\varepsilon_1, \varepsilon_2 \ldots \varepsilon_m$ are all UNOBSERVABLE random variables

Three Assumptions of factor model are:
- $E[F] = 0$, $Cov[F] = I = \begin{pmatrix} 1 & \cdots & 0 \\ 0 & \cdots & 1 \end{pmatrix}$
- $E[\varepsilon] = 0$, $Cov[\varepsilon] = \psi = \begin{pmatrix} \psi_1 & 0 & 0 \\ 0 & \psi_2 & \\ 0 & 0 & \psi_3 \end{pmatrix}$
- $f$ and $\varepsilon$ are independent, $Cov[F, \varepsilon] = 0$

$\Rightarrow$ The above is the orthogonal factor model. If we assume $Cov[F]$ is not diagonal we have "OBLIQUE" factor model

- We have $\Sigma = Cov[X] = E[(X-\mu)(X-\mu)^T]$

$= E[(\Lambda f + \varepsilon)(\Lambda f + \varepsilon)^T]$

$= E[(\Lambda F)(\Lambda F)^T + \varepsilon(\Lambda F)^T + (\Lambda F)\varepsilon^T + \varepsilon\varepsilon^T]$

$= \Lambda E[FF^T]\Lambda^T + E[\varepsilon F^T]\Lambda^T + \Lambda E[F\varepsilon^T] + E[\varepsilon\varepsilon^T]$

$= \Lambda\Lambda^T + \psi$

We can show $Cov[X, F] = E[(X-\mu)(F-0)] = \Lambda$, hence the covariance of the observed variable $X_i$ and the unobserved factor $F_j$ is the factor loading $\lambda_{ij}$.

VARIENCE
- Variance of $X$ can be split into two parts
- First portion of the variance for the $i^{th}$ component arisd from the $m$ common factors and is referred to as the $i^{th}$ communality
- Remainder of variance for the $i^{th}$ component is due to the specific factor referred to as the uniqueness

$\sigma_i^2 = \lambda_{i1}^2 \lambda_{i2}^2 + \ldots \lambda_{ip}^2 + \psi_i$

$= h_i^2 + \psi_i$

denoting $i^{th}$ communality as $h_i^2$

$Var[X_i] = $ communality $+$ uniqueness

- The $i^{th}$ communality is the sum of square of the loading of the $i^{th}$ variable on the $p$ common factors

MLA

FACTOR ANALYSIS

- Factor model assumes that the $m + m(m-1)/2$ variances and covariances of $X$ can be reproduced from the $mp$ factor loadings $\lambda_{ij}$ and the $m$ specific variances $\psi_i$.

- In the case where $p \ll m$ the factor model provides a simplified version of the covariance in $X$ with fewer parameters than the $m(m+1)/2$ parameters in $\Sigma$.

- Example: if $X = (X_1 \ldots X_{12})$ and a factor model with $p=2$ is appropriate, then the $12 \times 13/2 = 78$ elements of $\Sigma$ are described in terms of the $pm + m = 2 \times 12 + 12 = 36$ parameters $\lambda_{ij}$ and $\psi_i$ of the factor model.

- Unfortunately most covariance matrices cannot be uniquely factored as $\Lambda\Lambda' + \Psi$ where $p \ll m$.

- FA is not affected by re-scaling of the variables. Rescaling $X$ is the equivalent to letting $y = CX$ where $C = \text{diag}(C_i)$.

If the factor model holds with $\Lambda = \Lambda_X$ and $\Psi = \Psi_X$

$$y - C\mu = C\Lambda_X f + C\varepsilon$$

So, $\text{Var}[y] = C\Sigma C = C\Lambda_X \Lambda_X^\top C + C\Psi_X C$

Adding a $C$ in everywhere

- To demonstrate that most covariance matrices $\Sigma$ cannot be uniquely factored as $\Lambda\Lambda' + \Psi$ where $p \ll m$, let $G$ be any $m \times m$ orthogonal matrix $(GG' = I)$. Then

$$X - \mu = \Lambda f + \varepsilon$$
$$= \Lambda GG^\top f + \varepsilon$$
$$= \Lambda^* f^* + \varepsilon$$

- Here $\Lambda^* = \Lambda G$ and $f^* = G^\top f$

- It follows that $E[f^*] = 0$ and $\text{cov}[f^*] = I$

- Thus it is impossible, given the data $X$, to distinguish between $\Lambda$ and $\Lambda^*$ since they both generate the same cov matrix $\Sigma$: $\Sigma = \Lambda\Lambda^\top + \Psi$
$$= \Lambda GG^\top \Lambda' + \Psi = \Lambda^* \Lambda^{*\top} + \Psi$$

- This leads to idea of factor rotation, since orthogonal matrices correspond to rotations of the coordinate system of $X$

- Usually we estimate one possibility for $\Lambda$ and $\Psi$ and rotate the resulting loading matrix (multiply by an orthogonal matrix $G$) so as to ease interpretation

- Once the loadings and specific variances are estimated, estimated values for the factors themselves (called factor scores) are constructed

- When the data $X$ is assumed to be normally distributed, then estimates for $\Lambda$ and $\Psi$ can be obtained by MLE

- Normal location parameter is replaced by its MLE $\bar{X}$, whilst the log-likelihood of the data depend on $\Lambda$ and $\Psi$ through $\Sigma$

- To ensure $\Lambda$ is well defined (invariance is caused by orthogonal transformations) the computationally convenient (unly verted) condition is enforced

$$\Lambda^T \Psi^{-1} \Lambda = \Delta \quad \text{where } \Delta \text{ is a diagonal matrix}$$

- Numerical optimisation of log likelihood performed to obtain MLE $\hat{\Lambda}$ and $\hat{\Psi}$

- Problems occur if $\Psi_{ii} = 0$ (when uniqueness is so) the variance in the variable $X_i$ is completely accounted for by the factor $f_i$.

- problem if $\Psi_{ii} < 0 \rightarrow$ Heywood Case

$\Rightarrow$ occurs if there are too many common factors, too few common factors, not enough data or inappropriate application of model for data

$\Rightarrow$ example resolved (this) by resetting $\Psi_{ii}$ to be a small positive number.

## Factor Rotation

- If the initial loadings are subject to an orthogonal transformation (i.e. multiply by orthogonal matrix $G$) the covariance can still be reproduced

- An orthogonal transformation corresponds to a rigid rotation or reflection of coord axis.

- Hence orthogonal transformation of factor loadings known as factor rotation

- Doesn't matter if we use $\Lambda$ or $\Lambda^x = \Lambda G$ Since

$$\Sigma = \Lambda \Lambda^T + \Psi = \Lambda G G^T \Lambda^T + \Psi = \Lambda^x \Lambda^{xT} + \Psi$$

- One rotation may be more useful than another

MLA $\overset{5}{}$

# FACTOR ANALYSIS

## Controversial?
- Some say we are manipulating results, others say sharpening focus. Can be used as an intermediate step to reduce dimensionality prior to other techniques.

## Simple Structure
- One idea - rotate the factors so that each variable has a large loading on a single factor and small loadings on the others. Variables can then be split into disjoint sets, each associated with one factor.
- A factor $j$ can be interpreted as an average quantity over those variables $i$ where $\lambda_{ij}$ is large.

## Kaiser's Varimax rotation
- Squared loading $\lambda_{ij}^2$ is the proportion of the variance in variable $i$ that is attributed to common factor $j$.

$$\text{Var}[x_i] = \lambda_{i1}^2 + \lambda_{i2}^2 + \ldots \lambda_{ip}^2 + \psi_i$$
$$= h_i^2 + \psi_i$$

- Aim for a rotation that makes the squared loading $\lambda_{ij}^2$ either large or small, i.e. not medium sized values.
- Let $\tilde{\lambda}_{ij}^x = \lambda_{ij}^x / h_i$ be the final rotated loadings scaled by $\sqrt{}$ of communality.

Varimax procedure selects the orthogonal transformation $G$ maximising the sum of the column variances across all factors $j = 1 \ldots p$

$$V = \frac{1}{m} \sum_{j=1}^{p} \sum_{i=1}^{m} \tilde{\lambda}_{ij}^{x4} - \frac{1}{m^2} \sum_{j=1}^{p} \left( \sum_{i=1}^{m} \tilde{\lambda}_{ij}^x \right)^2$$

- Maxing $V$ corresponds to "spreading out" the squares of the loadings on each factor as much as possible so that both groups of large and negligible coefficients are found in any column of $\Lambda^x$
- Scaling the rotated loadings has the effect of giving variables with small communalities relatively more weight. After $G$ has been determined the loadings $\tilde{\lambda}_{ij}^x$ are multiplied by $h_i$ to ensure original communalities are preserved.

## OBLIQUE ROTATIONS

- An oblique rotation responds to a non rigid rotation of the coordinate system
  ie resulting axes need no longer be perpendicular
- oblique relaxed the orthogonally constraint in order to gain simplicity of interpretation
- Example: promax rotation from procrustian rotation
- less popular than orthogonal

## Factor Scores

- Interest usually lie in the parameters of a factor model (ie. $\lambda_{ij}$ and $\varepsilon_i$) but the estimated values of the common factors (ie factor scores ?) may be required
- location of each original ob in reduced factor space is often necessary as input for a subsequent analysis
- One method to use to estimate factor scores is regression method
  For MVN setting can be shown $\hat{f} = \Lambda^T \Sigma^{-1} X$

## Interpretation

- communality as the sum of square loading in each row of loading matrix
- Communality for each variable added to the uniqueness for each variable is the total variance for that variable ($=1$ as $R$ standardised data)
- The $i^{th}$ loading value are the sum of squared loading for that factor
- Proportion of total standardised sample variance due to $j^{th}$ factor is:

$$\frac{\hat{\lambda}^2_{ij} + \hat{\lambda}^2_{2j} + \dots + \hat{\lambda}^2_{ms}}{\sigma_{11} \sigma_{22} + \dots + \sigma_{mm}} = \frac{\hat{\lambda}^2_{ij} \dots + \hat{\lambda}^2_{ms}}{m} \qquad m = 10 \; \text{etc.}$$

## PCA v Factor Analysis

- PCA looks for linear combinations of data matrix $X$ that are uncorrelated and high variance, FA seeks unobserved linear combinations of variables representing underlying fundamental quantities
- PCA makes no assumption about form of covariance matrix, FA assumes data comes from a well defined model in which specific assumptions hold eg. $E[F] = 0$
- PCA: data $\Rightarrow$ PC's      FA: factors $\Rightarrow$ data
- When specific variances are large, they are absorbed into PC's, FA makes special provision for them. When specific variances are small, PCA and FA give similar results
- two analysis often performed together. Example could use PCA to determine number of factors to extract in FA study

MLA

08/05/15  FACTOR ANALYSIS

- Mathematical approach for attempting to explain the correlation between a large set of variables in terms of a small number of underlying factors
- Main assumption of FA is that it is not possible to observe these underlying factors directly, i.e. they are 'latent'.

- Dim reduction tech, share some obj a) PCA
- Attempt to denote the relationships in a large set of variables through a smaller number of dimensions
- FA much more elaborate

Example: performance a group of children in Irish, eng and math) record correlation matrix for 3 (ue) vars:

$$\begin{pmatrix} 1 & 0.83 & 0.78 \\ & 1 & 0.67 \\ & & 1 \end{pmatrix}$$

- Dimensionality of this matrix can be reduced from m=3 to m=1 by expressing variables $\lambda$

$$X_1 = \lambda_1 F + \varepsilon_1$$
$$X_2 = \lambda_2 F + \varepsilon_2$$
$$X_3 = \lambda_3 F + \varepsilon_3$$

F is the underlying common factor
$\lambda_i$'s are known as factor loadings
$\varepsilon_i$'s are errors or specific factors
- common factor will be given an interpretation e.g general ability
- Specific factor $\varepsilon_i$ will have small value if $X_i$ is closely related to general ability

The Factor Model
- The observable random vector $X = (X_1 ... X_m)$ mean $\mu$ cov matrix $\Sigma$
- FM states that X is linearly dependent upon a few unobservable random variables $f_1, f_2 ... f_p$ called common factors and m additional sources of variation $\varepsilon_1, \varepsilon_2 ... \varepsilon_m$ called specific factors

$$X_i - \mu_i = \lambda_{i1} f_1 + \lambda_{i2} f_2 + ... + \lambda_{ip} f_p + \varepsilon_i$$

$$X - \mu = \Lambda F + \varepsilon$$

$\lambda_{ij}$ is factor loading of $i^{th}$ variable on the $j^{th}$ factor.
$\Lambda$ = matrix of factor loadings
$i^{th}$ specific factor $\varepsilon_i$ is associated only with repeat $X_i$
Note that $f_1, f_2 ... f_m, \varepsilon_1 ... \varepsilon_m$ are all unobserved random variables

Under FM assume:

$$\mathbb{E}[f] = 0 \qquad Cov[f] = I = \begin{pmatrix} 1 & 0 & 0 \\ & 1 & \\ 0 & 0 & 1 \end{pmatrix}$$

2. $\mathbb{E}[\varepsilon] = 0$    $Cov[\varepsilon] = \Psi = \begin{pmatrix} \Psi_1 & 0 & 0 & 0 \\ 0 & \Psi_2 & 0 & 0 \\ 0 & 0 & \ddots & \\ 0 & 0 & & \Psi_p \end{pmatrix}$

3. f and $\varepsilon$ are independent, $Cov[f, \varepsilon] = 0$.

Above 1 referred to as orthogonal factor model. However if it is assumed that $Cov[f]$ is not diagonal, then the model is referred to as the oblique factor model

- Orthogonal FM implies a specific cov structure for $X$: $= \Lambda \Lambda^T + \Psi$
- Can also be shown that $Cov[X, F] = \mathbb{E}[(x-\mu)(f-0)] = \Lambda$ hence the cov of the observed variable $X_i$ and unobserved factor $f_j$ is the factor loading $\lambda_{ij}$.

Variance
- Variance can be split in 2.
- First part for $i^{th}$ component arises from the m common factors referred to as the $i^{th}$ communality
- The remainder of variance is due to the specific factor referred to as uniqueness
- Denote $i^{th}$ communality by $h_i^2$ etc:
  $$\sigma_i^2 = \lambda_{i1}^2 + \lambda_{i2}^2 + \dots \lambda_{ip}^2 + \Psi_i$$
  $$= h_i^2 + \Psi_i$$
- $i^{th}$ communality is the sum of squares of loading of the $i^{th}$ variable onto the p common factors

- FM assumes that the $m+m(m-1)/2$ variance and covariance of $X$ can be reproduced from the mp factor loadings $\lambda_{ij}$ and the m specific variance $\Psi_i$
- In case when $p < m$ the FM provides a simpler way of the covariance in $X$ with fewer parameters than the $m(m+1)/2$ parameters in $\Sigma$

- Example if $X = (X_1, \dots X_{12})$ FM with $p = 2$ is used. Then the $12 \times 13/2 = 78$ elements of $\Sigma$ are described in term of the $pm+m = 2 \times 12 + 12 = 36$ parameters $\lambda_{ij}$ and $\Psi_i$ of the factor model

- Most cov matrix cant be uniquely factored as $\Lambda \Lambda^T + \Psi$ with $p < m$

Scale Invariance
- Rescaling variable of $X$ is equivalent to letting $Y = CX$ where $C = diag(c_i)$.
- FA not affected by re-scaling of variables
- leads to idea of factor rotation

ML for FA
- When data $X$ is assumed to be normally distributed, then estimates for $\Lambda$ and $\Psi$ can be obtained by maximum likelihood
- Uniqueness condition entail $\Lambda^T \Psi^{-1} \Lambda = \Delta$

# MLA

09/05/15 **FACTOR ANALYSIS**

- Problems may occur if $\psi_{ii} = 0$
- This happens when uniqueness is zero i.e. variance in the variable $x_i$ is completely accounted for by the factor $f_i$.
- May find $\psi_{ii} < 0$ improper solution called Haywood case
- Occurs if there are too many common factors, too few common factors, not enough data or the inappropriate application or the model for the data.
- Resolved by re-setting $\psi_{ii}$ to a small positive number before proceeding

## Factor Rotation

- If initial loadings are subject to an orthogonal transformation (i.e. multiplied by an orthogonal matrix $G$), the cov matrix $\Sigma$ can still be reproduced
- An orthogonal transformation corresponds to a rigid rotation or reflection of coordinate axes
- Called factor rotation
- It is immaterial whether $\Lambda$ or $\Lambda^* = \Lambda G$ is reported since:
  $$\Sigma = \Lambda \Lambda^T + \psi = \Lambda G G^T \Lambda^T + \psi = \Lambda^* \Lambda^{*T} + \psi$$
- Argued that this is manipulating result
- Others say it is sharpening the focus
- FA generally used as an intermediate step to reduce the dimensionality of a data set prior to clear statistical analysis

## A Simple structure

- One idea would be to rotate the factors so that each variable has a large loading on a single factor and small loadings on the others.
- Variables can then be split into disjoint sets, each of which is associated with one factor.
- A factor can then be interpreted as an clearly quality of those variables for which $\lambda_{ij}$ is large

## Kaiser's Varimax Rotation

- Squared loading $\lambda_{ij}^2$ is the proportion of variance in variable $i$ that is attributable to common factor $j$:
  $$\text{var}[x_i] = \lambda_{i1}^2 + \lambda_{i2}^2 + \ldots + \lambda_{ip}^2 + \psi_i$$
  $$= h_i^2 + \psi_i$$
- We aim for a rotation that makes squared loadings $\lambda_{ij}^2$ either large or small i.e. few medium size values
- Let $x_{ij} = \lambda^*_{ij}/h_i$ be the final rotated loadings scaled by the square root of the communalities
- Varimax procedure selects the orthogonal transformation $G$ maximising the sum of the column variances across all factors $j=1 \ldots p$
- Maximising $V$ corresponds to spreading out the squared of the loadings on each factor or MAD
- Scaling rotated loadings has the effect of giving variables with small communalities relatively more weight. After $G$ has been determined the loadings $\lambda^*_{ij}$ are multiplied by $h_i$ to entry original communalities are preserved

## Oblique Rotation

- Degree of correlation allowed between factors is generally small (2 highly correlated factors = one factor)
- - results and need no longer be perpendicular.
- OR therefor relax the orthogonality constraint in order to gain simplicity in interpretation
- Promax rotation ⇒ derived from procrustes rotation

- Communalities are the sum of squared loadings in each row of loadings matrix
- Communality for each variable added to the uniqueness for each variable is total variance for that variable = 1 because of standardise
- Σ loadings cols are the sum of squared loadings for that factor
- Proportion variance = Σ loadings / m    m=10 in eg.
- Interpret like PCA.

## PCA vs Factor Analysis

- Both PCA and FA have similar aims
- PCA looks for linear combination of the data matrix X that are uncorrelated and of high variance, while FA seeks unobserved linear combination of the variables representing underlying fundamental quantities.

- PCA makes no assumption about the form of cov mat. FA assumes that the data comes from a well defined model in which specific assumptions hold.
  $E[F] = 0$ etc.

- PCA: data ⇒ PCs        FA: factors ⇒ data

- When specific variance are large they are absorbed into the PCs whereas FA make special provision for them. When specific variances are small, PCA and FA give similar results

- 2 analysis often performed together. Eg. PCA first to determine number of factors to extract in a FA study.

| | |
|---|---|
| - one linear Regression | Metric/Non metric |
| - Assumptions | Classical/(Least Square)   Kruskal |
| - Variance | Stress    Sammon stress |
| - Orthogonal Varimax Rotation | Procrustes |
| - Oblique | PCA vs MDS. |
| - PCA vs FA. | |