

Multivariate Statistical Analysis

Lecture 16

Fudan University

luoluo@fudan.edu.cn

1 Factor Analysis

2 Probabilistic Principle Component Analysis

1 Factor Analysis

2 Probabilistic Principle Component Analysis

Factor Analysis

Let the observable vector $\mathbf{y} \in \mathbb{R}^p$ be written as

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \boldsymbol{\mu} + \boldsymbol{\epsilon},$$

where

- ① $\mathbf{W} \in \mathbb{R}^{p \times q}$ is the loading matrix (parameter),
- ② $\mathbf{x} \in \mathbb{R}^q$ is the common factor (parameter/random vector),
- ③ $\boldsymbol{\mu} \in \mathbb{R}^p$ is the mean vector (parameter),
- ④ $\boldsymbol{\epsilon} \in \mathbb{R}^p$ is the specific factor (random vector).

The model is similar to regression, but \mathbf{x} is unobserved.

Factor Analysis

Example of sports games:

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \mu + \epsilon.$$

- ① \mathbf{y} : performance in real-world
- ② \mathbf{W} : system of the game
- ③ \mathbf{x} : attributes in the game
- ④ μ : average attributes
- ⑤ ϵ : noise/exception

Yao Ming's Stats

57 Games	82 Games	80 Games	48 Games
22.3 Points	17.5 Points	18.3 Points	25.0 Points
51.9% FG%	52.2% FG%	55.2% FG%	51.6% FG%
0.0% 3PT%	0.0% 3PT%	0.0% 3PT%	0.0% 3PT%
85.3% FT%	80.9% FT%	86.2% FT%	86.2% FT%
10.2 Rebounds	9.4 Rebounds	9.4 Rebounds	9.4 Rebounds
1.5 Assists	1.5 Assists	2.0 Assists	2.0 Assists
0.5 Steals	0.5 Steals	0.4 Steals	0.4 Steals
1.6 Blocks	2.0 Blocks	2.0 Blocks	2.0 Blocks
2.6 Turnovers	2.5 Turnovers	3.5 Turnovers	3.5 Turnovers
25.6 PER	21.9 PER	26.5 PER	26.5 PER
0.211 WS/48	0.202 WS/48	0.220 WS/48	0.220 WS/48



Example of mental tests for $\mathbf{y} = \mathbf{W}\mathbf{x} + \boldsymbol{\mu} + \epsilon$:

- ① Each component of \mathbf{y} is a (centralized) score on a battery of tests.
- ② The components of \mathbf{x} are the scores of the mental factors, linear combinations of these enter into the test scores.
- ③ Each component of $\boldsymbol{\mu}$ is the average score in the population.
- ④ The coefficients of these linear combinations are the elements of \mathbf{W} , and these are called factor loadings (common factors).
- ⑤ A component of ϵ is the part of the test score not “explained” by the common factors (error).

Example of recommending system for $\mathbf{y} = \mathbf{W}\mathbf{x} + \boldsymbol{\mu} + \epsilon$:

- ① Each component of \mathbf{y} is a (centralized) score on an item.
- ② The components of \mathbf{x} are the attributes of the user.
- ③ Each component of $\boldsymbol{\mu}$ is the average score in the population.
- ④ The coefficients of these linear combinations are the elements of \mathbf{W} , and these are called factor loadings (common factors).
- ⑤ The components of ϵ are noise.

Factor Analysis

The columns of $\mathbf{W} \in \mathbb{R}^{d \times q}$ establish an q -dimensional subspace of \mathbb{R}^d .

- ① This subspace is called the factor space.
- ② Vector $\mathbf{x} \in \mathbb{R}^q$ can be viewed as coordinates of a point in factor space.

1 Factor Analysis

2 Probabilistic Principle Component Analysis

Probabilistic Principle Component Analysis

Let $\mathbf{y}_1, \dots, \mathbf{y}_N$ be N independent observations and we have

$$\mathbf{y}_\alpha = \mathbf{W}\mathbf{x}_\alpha + \boldsymbol{\mu} + \epsilon_\alpha,$$

where $\mathbf{x}_\alpha \sim \mathcal{N}_q(\mathbf{0}, \mathbf{I})$ and $\epsilon_\alpha \sim \mathcal{N}_p(\mathbf{0}, \sigma^2 \mathbf{I})$ are independent for some $\sigma^2 > 0$.

We target to estimate parameters

$$\mathbf{W} \in \mathbb{R}^{p \times q}, \quad \boldsymbol{\mu} \in \mathbb{R}^p \quad \text{and} \quad \sigma^2 \in \mathbb{R}$$

by maximum likelihood estimation.

We are interested in the case of $q < p$.

Probabilistic Principle Component Analysis

Then, we have $\mathbf{y}_\alpha \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$, where $\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I}$.

The log-likelihood function is

$$-\frac{Nd \ln(2\pi)}{2} - N \ln \det(\mathbf{C}) - \text{tr}\left(\mathbf{C}^{-1} \sum_{\alpha=1}^N (\mathbf{y}_\alpha - \boldsymbol{\mu})(\mathbf{y}_\alpha - \boldsymbol{\mu})^\top\right).$$

The Maximum Likelihood Estimators

The maximum likelihood estimators of $\boldsymbol{\mu}$, \mathbf{W} and σ^2 are

$$\boldsymbol{\mu} = \bar{\mathbf{y}} = \frac{1}{N} \sum_{\alpha=1}^N \mathbf{y}_{\alpha}, \quad \hat{\mathbf{W}} = \mathbf{U}_q (\boldsymbol{\Lambda}_q - \hat{\sigma}^2 \mathbf{I}) \mathbf{R} \quad \text{and} \quad \hat{\sigma}^2 = \frac{1}{d-q} \sum_{j=q+1}^d \lambda_j,$$

where $\mathbf{U}_q \in \mathbb{R}^{d \times q}$ with columns are the principal eigenvectors of

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{N} \sum_{\alpha=1}^N (\mathbf{y}_{\alpha} - \bar{\mathbf{y}})(\mathbf{y}_{\alpha} - \bar{\mathbf{y}})^{\top},$$

$\boldsymbol{\Lambda}_q \in \mathbb{R}^{q \times q}$ is diagonal matrix with corresponding eigenvalues $\lambda_1, \dots, \lambda_q$ and \mathbf{R} is any $q \times q$ orthogonal matrix.

The Maximum Likelihood Estimators

The MLE estimator also minimize the Frobenius norm error

$$(\hat{\mathbf{W}}, \hat{\sigma}^2) = \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times q}, \sigma^2 \in \mathbb{R}^+} \left\| \hat{\mathbf{\Sigma}} - (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}) \right\|_F.$$

Lemma 1

Let $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}$ and $q = \min\{m, n\}$. Define the diagonal matrix $\mathbf{\Sigma}(\mathbf{A})$ whose (i, i) -th element is the i -th singular value of \mathbf{A} and the others are zero. We define $\mathbf{\Sigma}(\mathbf{A})$. Then we have

$$\|\mathbf{A} - \mathbf{B}\| \geq \|\mathbf{\Sigma}(\mathbf{A}) - \mathbf{\Sigma}(\mathbf{B})\|.$$

for every unitarily invariant norm.

The EM Algorithm

For the model

$$\mathbf{y} = \mathbf{W}\mathbf{x} + \boldsymbol{\mu} + \epsilon,$$

where $\mathbf{x} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{I})$ and $\epsilon \sim \mathcal{N}_d(\mathbf{0}, \sigma^2 \mathbf{I})$ are independent.

View $\{\mathbf{x}_\alpha\}_{\alpha=1}^N$ as missing data and $\{\mathbf{x}_\alpha, \mathbf{y}_\alpha\}_{\alpha=1}^N$ as the complete data.

- ① $\mathbf{y} \mid \mathbf{x} \sim \mathcal{N}_d(\mathbf{W}\mathbf{x} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$
- ② $\mathbf{x} \mid \mathbf{y} \sim \mathcal{N}_q(\mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{y} - \boldsymbol{\mu}), \sigma^2 \mathbf{M}^{-1})$, where $\mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}$

The EM Algorithm

The update of the EM algorithm

- 1 In E-step, we take the expectation

$$l_C = \mathbb{E} \left[\ln \left(\prod_{\alpha=1}^N p(\mathbf{x}_\alpha | \mathbf{y}_\alpha) \right) \right].$$

- 2 In the M-step, we maximized l_C with respect to \mathbf{W} and σ^2 :

$$\begin{aligned} \tilde{\mathbf{W}} &= \hat{\Sigma} \mathbf{W} (\sigma^2 \mathbf{I} + \mathbf{M}^{-1} \mathbf{W}^\top \hat{\Sigma} \mathbf{W})^{-1}, \\ \tilde{\sigma}^2 &= \frac{1}{d} \text{tr} \left(\hat{\Sigma} - \hat{\Sigma} \mathbf{W} \mathbf{M}^{-1} \tilde{\mathbf{W}}^\top \right). \end{aligned}$$

Note that the computational complexity of EM is $\mathcal{O}(Ndq)$, while MLE requires $\mathcal{O}(Nd^2 + d^3)$.