

Mixtures of unrestricted skew normal factor analyzers with missing information 县遺失訊息之混合非限制型的偏斜常態因子分析器 Student: Hong-Ying Chen Advisor: Tsung-I Lin

Institute of Statistics, National Chung Hsing University, Taichung, Taiwan

國立中興大學



Abstract

Mixtures of factor analyzers (MFA) based on the restricted skew normal distribution (rMSN) has been shown to be a flexible tool for modeling asymmetrical high-dimensional data with heterogeneity. However, the rMSN distribution is oft-criticized a lack of sufficient ability to accommodating skewness arising from more than one feature space. This thesis presents an alternative extension of MFA by assuming the unrestricted skew normal (uMSN) distribution for the component factors. In particular, the proposed mixtures of uMSN factor analyzers (MuSNFA) can simultaneously accommodate multiple directions of skewness and the occurrence of missing values or nonresponses. Under the missing at random mechanism, we develop a computationally feasible expectation conditional maximization (ECM) algorithm for computing the maximum likelihood estimates of model parameters. Practical aspects related to model-based clustering, prediction of factor scores and missing values are also discussed. The utility of the proposed methodology is illustrated with the analysis of simulated data and an application to the Pima Indian women diabetes data containing genuine missing values.

Motivation

Unlike the rMSN distribution, the uMSN distribution is able to capture skewness in more than one direction.

Aim: propose a skew extension of MFA based on the uMSN distribution and allows the handling of missing values.



Methodology

> MuSNFA model

 $Y_i = \mu_i + B_i U_{ii} + \varepsilon_{ii}$, with probability π_i (i = 1, ..., g) $\boldsymbol{U}_{ij} \sim uSN_q \left(-c\boldsymbol{\Delta}_i^{-1/2}\boldsymbol{\Lambda}_i \boldsymbol{1}_q, \boldsymbol{\Delta}_i^{-1}, \boldsymbol{\Delta}_i^{-1/2}\boldsymbol{\Lambda}_i \right), \ \varepsilon_{ij} \sim N_p(\boldsymbol{0}, \boldsymbol{D}_i), \ \boldsymbol{U}_{ij} \perp \varepsilon_{ij}$ where Λ_i be a q^*q skewness parameter matrix, and $\Delta_i = I_q + (1 - c^2)\Lambda_i\Lambda_i^T$ with $c = \sqrt{2/\pi}$. rightarrow Introduce two permutation matrices $\mathbf{O}_{j}(p_{j}^{o} \times p)$ and $M_{j}((p - p_{j}^{o}) \times p)$ such that $Y_{i} = O_{i}^{T} Y_{i}^{o} + M_{i}^{T} Y_{i}^{m}$. The hierarchical representation of MuSNFA is

 $Y_j^o|(Z_{ij}=1) \sim CFUSN_{p_i^o,q}(\mu_{ij}^o - c\alpha_{ij}^o \mathbf{1}_q, \Sigma_{ij}^{oo}, \alpha_{ij}^o)$ $\mathbf{Z}_i \sim \mathcal{M}(1; \pi_1, \dots, \pi_q),$ where $Y_i^o = O_j Y_j, Y_j^m = M_j Y_j, \ \mu_{ij}^o = O_j \mu_i, \ \Sigma_{ij}^{oo} = O_j \Sigma_i O_j^T, \ \alpha_{ij}^o = O_j \alpha_i, \ \Sigma_i = B \Delta_i^{-1} B^T + D_i$ $\alpha_i = B \Delta_i^{-1/2} \Lambda_i$, and Z_i is a set of binary allocation indicators. rightarrow The component pdf of $Y_i^o | (Z_{ij} = 1)$ is given by $\boldsymbol{\phi}(\boldsymbol{y}_{j}^{o};\boldsymbol{\theta}_{i}) = 2^{q} \boldsymbol{\phi}_{p_{i}^{o}}\left(\boldsymbol{y}_{j}^{o};\boldsymbol{\mu}_{ij}^{o} - \boldsymbol{c}\boldsymbol{\alpha}_{ij}^{o}\boldsymbol{1}_{q},\boldsymbol{\Omega}_{ij}^{oo}\right) \times \boldsymbol{\Phi}_{q}(\boldsymbol{\alpha}_{ij}^{o^{\top}}\boldsymbol{\Omega}_{ij}^{oo}^{-1}(\boldsymbol{y}_{j}^{o} - \boldsymbol{\mu}_{ij}^{o} + \boldsymbol{c}\boldsymbol{\alpha}_{ij}^{o}\boldsymbol{1}_{q});\boldsymbol{I}_{q} - \boldsymbol{\alpha}_{ij}^{o^{\top}}\boldsymbol{\Omega}_{ij}^{oo^{-1}}\boldsymbol{\alpha}_{ij}^{o}),$ where $\Omega_{ij}^{oo} = O_j \Omega_i O_j^T$ with $\Omega_i = \Sigma_i + \alpha_i \alpha_i^T$. For ease of estimation, we reparametrize

$$\widetilde{B_i} \triangleq B_i \Delta_i^{-1/2}$$
 and $\widetilde{U}_{ij} \triangleq \Delta_i^{1/2} U_{ij}$





Simulation & Real data

> Simulation

- A simulation study is undertaken to examine the performance of the proposed MuSNFA model in comparison with the MFA and MrSNFA approaches under varying proportions of synthetic missing values.
- 100 Monte Carlo samples are created from the 3-component MFA model with p=5 and q=2 factors, while the latent factors are independently generated from the chi-square distribution with one degree of freedom χ_1^2 for yielding strong effects of skewness and kurtosis. MuSNFA model outperforms the other two competing approaches in almost all trials. The MFA model has the worst performance due to a lack of sufficient robustness against serious departure of normality assumption.



| Missing | Criterion | | | n = 300 | | | n = 600 | | | n = 1200 | |
|---------|-----------|------|---------|---------|---------|---------|---------|---------|----------|----------|----------|
| rate | | | MFA | MrSNFA | MuSNFA | MFA | MrSNFA | MuSNFA | MFA | MrSNFA | MuSNFA |
| r = 10% | BIC | Mean | 4790.81 | 4756.70 | 4735.59 | 8958.33 | 8828.85 | 8782.34 | 17664.59 | 17352.88 | 17243.15 |
| | | Freq | 4 | 14 | 82 | 0 | 4 | 96 | 0 | 6 | 94 |
| | | Mean | 4902.15 | 4864.91 | 4834.12 | 9246.50 | 9112.59 | 9022.93 | 18265.53 | 17944.62 | 17717.26 |
| | ICL | Freq | 6 | 11 | 83 | 0 | 6 | 94 | 0 | 2 | 98 |
| | | Mean | 0.834 | 0.842 | 0.870 | 0.834 | 0.858 | 0.886 | 0.832 | 0.873 | 0.911 |
| | CCR | Freq | 25 | 19 | 56 | 12 | 25 | 63 | 8 | 5 | 87 |
| | | Mean | 0.623 | 0.634 | 0.689 | 0.617 | 0.646 | 0.710 | 0.617 | 0.675 | 0.759 |
| | ARI | Freq | 23 | 17 | 60 | 12 | 24 | 64 | 8 | 3 | 89 |
| | | Mean | 1.184 | 1.139 | 1.112 | 1.183 | 1.151 | 1.130 | 1.142 | 1.097 | 1.066 |
| | MSE | Freq | 14 | 35 | 51 | 6 | 28 | 66 | 2 | 16 | 82 |
| r = 30% | | Mean | 3757.54 | 3702.84 | 3693.55 | 7281.53 | 7134.67 | 7113.46 | 14331.74 | 13997.08 | 13938.48 |
| | BIC | Freq | 0 | 16 | 84 | 0 | 12 | 88 | 0 | 9 | 91 |
| | | Mean | 3972.24 | 3882.21 | 3871.68 | 7747.61 | 7552.67 | 7513.12 | 15312.47 | 14892.26 | 14750.88 |
| | ICL | Freq | 0 | 24 | 76 | 0 | 20 | 80 | 0 | 7 | 93 |
| | | Mean | 0.781 | 0.745 | 0.766 | 0.803 | 0.759 | 0.803 | 0.813 | 0.782 | 0.826 |
| | CCR | Freq | 49 | 16 | 35 | 34 | 5 | 61 | 28 | 12 | 60 |
| | | Mean | 0.475 | 0.439 | 0.473 | 0.508 | 0.451 | 0.523 | 0.526 | 0.486 | 0.563 |
| | ARI | Freq | 47 | 16 | 37 | 32 | 5 | 63 | 25 | 14 | 61 |
| | | Mean | 1.529 | 1.485 | 1.457 | 1.432 | 1.393 | 1.379 | 1.419 | 1.353 | 1.335 |
| | MSE | Freq | 16 | 29 | 55 | 15 | 26 | 59 | 1 | 28 | 71 |

It follows that $Y_i = \mu_i + \widetilde{B}_i \widetilde{U}_{ii} + \varepsilon_{ii}$, where $\widetilde{U}_{ii} \sim uSN_a(-c\Lambda_i \mathbf{1}_a, I_a, \Lambda_i)$.

 $\mathbf{\Theta} = \{\pi_1, \dots, \pi_{g-1}, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_g\}$ represent all the unknown parameters of the mixture model, where $\boldsymbol{\theta}_{i} = (\boldsymbol{\mu}_{i}, \boldsymbol{B}_{i}, \boldsymbol{D}_{i}, \boldsymbol{\Lambda}_{i})$.

The complete data likelihood function of
$$\Theta$$
 given $Y_c = (y^o, y^m, \widetilde{U}, \gamma, Z)$ is

$$L_c(\Theta \mid Y_c) = \prod_{\substack{j=1\\n}}^n f(y_j^o \mid Z_{ij} = 1) f(Z_j) = \prod_{\substack{j=1\\i=1}}^n \prod_{\substack{i=1\\i=1}}^g \{\pi_i \varphi(y_j^o \mid \theta_i)\}^{Z_{ij}}$$

$$= \prod_{\substack{j=1\\i=1}}^g \prod_{\substack{i=1\\i=1}}^g \{\pi_i \varphi_{p_j^o}(y_j^o; \mu_{ij}^o + \widetilde{B}_{ij}^o \widetilde{U}_{ij}, D_{ij}^{oo}) \varphi_q(\widetilde{U}_{ij}; -c\Lambda_i \mathbf{1}_q + \Lambda_i \gamma_j, I_q)\}^{Z_{ij}}$$

> **Proposition**:

For notational simplicity, the symbol "| \cdots " stands for conditioning on y_i^o and $Z_{ij} = 1$.



Pima Indians Diabetes Data Set (PIMA data)

 \checkmark There are (p=8) attributes measured on 500 non-diabetes and 278 diabetes female patients. We consider the fitting of 2-component MFA, MrSNFA and MuSNFA models to this dataset with *q* ranging from 1-4.

The factor scores predicted by the three approaches are all positively skewed, but the scattering patterns somewhat different. The MuSNFA offer the strongest degree of

| | Number of missing items | | | | | Patients | NA- | NA- | NA- | |
|--------------|-------------------------|---------|---------|--------|--------|----------|-----------|---------|--------|---------|
| | 0 | 1 | 2 | 3 | 4 | 5 | 1 attents | patient | var | obs |
| | 262 | 95 | 124 | 13 | 6 | 0 | 500 | 238 | 5 | 406 |
| Non-diabetes | (52.40) | (19.00) | (24.80) | (2.60) | (1.20) | (0.00) | (65.10) | (30.99) | (62.5) | (10.15) |
| Diabets | 130 | 47 | 75 | 15 | 1 | 0 | 268 | 138 | 5 | 246 |
| | (48.50) | (17.54) | (27.99) | (5.60) | (0.37) | (0.00) | (34.90) | (17.97) | (62.5) | (11.47) |
| The state | 392 | 142 | 199 | 28 | 7 | 0 | 768 | 376 | 8 | 652 |
| Total | (51.04) | (18.49) | (25.91) | (3.65) | (0.91) | (0.00) | (100) | (48.96) | (100) | (10.61) |

The percentages (%) are listed in parentheses (numbers/n for patient-wise; numbers/p for variable-wise; numbers/(np) for observation-wise

