Mixtures of factor-analytic models with covariates for modeling multiply censored dependent variables 具共變數之混合因子分析模型對於多重設限相依變數之建模 Student: Wan-Chen Hsieh Advisor: Tsung-I Lin Institute of Statistics, National Chung Hsing University, Taichung, Taiwan

Abstract

Censored data arise frequently in diverse applications in which observations to be measured may be subject to some upper and lower detection limits due to the restriction of experimental apparatus such that they are not exactly quantifiable. Mixtures of factor analyzers with censored data (MFAC) have been recently proposed for model-based density estimation and clustering of high-dimensional data under the presence of censored observations. In this paper, we consider an extended version of MFAC with covariates to accommodate multiply censored dependent variables and develop two analytically feasible EM-type algorithm for computing maximum likelihood estimates of the parameters with closed-form expressions. Moreover, we provide an information-based method to compute asymptotic standard errors of mixing proportions and regression coefficients. The utility and performance of the proposed methodology is illustrated through two real data examples.

Introduction

An extension of the MFAC approach (Wang et al. 2019) is introduced by including regression covariates. The proposed model clearly encompasses the MFAC approach with advantages including the flexibility in situations involving related independent variables within g classes. A description of finding ML estimates of model parameters through two EM-based algorithms in which the E-step relies on calculating the firsttwo moments of the truncated multivariate-normal distribution (Lin, 2009) is provided. The computational procedures we developed are rather stable and efficient for estimating the proposed model. In addition, we offer an approximation of the empirical Hessian matrix via the Louis' method (Louis, 1982), so that the standard errors of mixture regression coefficients and mixing proportions can be obtained. The utility of our methodology is illustrated through the analysis of two real-life datasets related to educational assessment and water quality.

Model Formulation

> MFA with linear regression model

 $m{y}_i = m{X}_im{eta}_i + m{B}_im{u}_{ij} + m{arepsilon}_{ij}$ with probability π_i (for *i*=1,...,*g*). $u_{ij} \sim N_q(\mathbf{0}, I_q)$ and $\varepsilon_{ij} \sim N_p(\mathbf{0}, D_i)$ are mutually independent for all *i* and *j*.

Censoring information

Define $v_j = (v_{j1}, ..., v_{jp})^{\dagger}$ as the vector of uncensored reading or censoring level, and $c_i = (c_{i1}, \dots, c_{ip})^{\dagger}$ be the vector of censoring indicators. We shall focus the case of left censoring:

 $y_{jk} \le v_{jk}, \quad \text{if} \quad c_{jk} = 1$ $y_{jk} = v_{jk}, \quad \text{if} \quad c_{jk} = 0$

Incorporating the censoring information into the model, we have the distribution, $y_j \mid (v_j, c_j, z_{ij} = 1) \sim TN_p(X_j \beta_i, \Sigma_i; A_j)$ where $TN_p(\mu, \Sigma; A_j)$ denotes a *p*-variate truncated normal distribution restricted within a right-truncated hyperplane $\mathbb{A}_j = \mathbb{A}_i^0 \times \mathbb{A}_i^c$, where $\mathbb{A}_{j}^{o} = \{-\infty < y_{jk} < \infty, k = 1, \dots, p_{j}^{o}\} \text{ and } \mathbb{A}_{j}^{c} = \{y_{jk} \le v_{jk}, k = 1, \dots, p - p_{j}^{o}\}.$ Therefore, the MFAC admits a two-level hierarchical representation:

 $\boldsymbol{v}_j \mid (\boldsymbol{c}_j, z_{ij} = 1) \sim f_{ij}(\boldsymbol{v}_j \mid \boldsymbol{c}_j, \boldsymbol{\theta}_i),$

 $\boldsymbol{z}_i \sim \mathcal{M}(1; \pi_1, \cdots, \pi_a).$ where $f_{ij}(\boldsymbol{v}_j \mid \boldsymbol{c}_j, \boldsymbol{\theta}_i) = f(\boldsymbol{y}_j^c \leq \boldsymbol{v}_j^c \mid \boldsymbol{y}_j^o = \boldsymbol{v}_j^o, \boldsymbol{\theta}_i) f(\boldsymbol{y}_j^o = \boldsymbol{v}_j^o \mid \boldsymbol{\theta}_i) = \Phi_{p-p_j^o}(\boldsymbol{v}_j^c \mid \boldsymbol{\mu}_{ij}^{c \cdot o}, \boldsymbol{\Sigma}_{ij}^{c \cdot o}) \phi_{p_j^o}(\boldsymbol{v}_j^o \mid \boldsymbol{\mu}_{ij}^o, \boldsymbol{\Sigma}_{ij}^{o o o})$ and $\mu_{ij}^{\circ} = O_j X_j \beta_i, \Sigma_{ij}^{\circ\circ} = O_j \Sigma_i O_j^{\top}, \mu_{ij}^{\circ\circ} = C_j X_j \beta_i + C_j \Sigma_i S_{ij}^{\circ\circ} (y_j - X_j \beta_i),$

 $S_{ij}^{\text{oo}} = O_j^{\top} (O_j \Sigma_i O_j^{\top})^{-1} O_j \text{ and } \Sigma_{ij}^{\text{cc·o}} = C_j (I_p - \Sigma_i S_{ij}^{\text{oo}}) \Sigma_i C_j^{\top}.$ > Given the hierarchical structures, we have the following conditional distributions $(\alpha, \alpha) = (\alpha, \alpha) + ($

Real Examples

EGRA data

We analyze the dataset presented by Costa et al. (2014) for 502 Peruvian students in 2007, who were taken by four of ten EGRA tasks, the four tasks in the dataset are:

(i) recognizing letters of the alphabet, (ii) recognizing simple words,

(iii) simple decoding meaningless words, and (iv) reading of a passage.

The design matrix X_i of dimension 4×4 for each student *j* corresponds to fixed effects vector $\beta_i = (\beta_{i1}, ..., \beta_{i4})$, for i = 1, ..., g.

The covariates in X_i related to β_i include gender (0=F, 1=M), grade (0=2nd, 1=3rd), residence **zone** (0 = Rural, 1 = Urban) and the adjusted **age** by subtracting the mean of students' ages. Censoring rate:

The first 10% of lowest total scores.



Water resources data

Our 2nd example concerns the water resources data (Hoffman and Johnson, 2015) that collected the trace metal concentration levels of certain dissolved trace metals in the freshwater. There are p=5 (trace metals) concentration levels from *n*=184 independent randomly selected sites in freshwater streams across Virginia.



• The figure shows the sequential points of five trace metals together with their detection limits and censoring proportions.

In this example, we are interested in comparing the convergence behavior of the EM and AECM algorithms for fitting the proposed model.

(a) g=1; q=1	(b) g=1; q=2

$$\begin{aligned} \mathbf{z}_{j} \mid (\mathbf{v}_{j}, \mathbf{c}_{j}) &\sim \mathcal{N}(1, \pi_{1j}, \dots, \pi_{gj}) \\ \mathbf{y}_{j}^{\mathrm{o}} \mid (z_{ij} = 1) &\sim N_{p_{j}^{\mathrm{o}}}(\boldsymbol{\mu}_{ij}^{\mathrm{o}}, \boldsymbol{\Sigma}_{ij}^{\mathrm{oo}}) \\ \mathbf{y}_{j}^{\mathrm{c}} \mid (\mathbf{y}_{j}^{\mathrm{o}}, z_{ij} = 1) &\sim N_{p-p_{j}^{\mathrm{o}}}(\boldsymbol{\mu}_{ij}^{\mathrm{co}}, \boldsymbol{\Sigma}_{ij}^{\mathrm{cco}}) \\ \mathbf{y}_{j}^{\mathrm{c}} \mid (\mathbf{y}_{j}^{\mathrm{o}}, \mathbf{v}_{j}, \mathbf{c}_{j}, z_{ij} = 1) &\sim TN_{p-p_{j}^{\mathrm{o}}}(\boldsymbol{\mu}_{ij}^{\mathrm{co}}, \boldsymbol{\Sigma}_{ij}^{\mathrm{cco}}; \mathbf{A}_{j}^{\mathrm{c}}) \\ \mathbf{u}_{ij} \mid (\mathbf{y}_{j}, z_{ij} = 1) &\sim N_{q}(\boldsymbol{\Gamma}_{i}^{\mathsf{T}}(\mathbf{y}_{j} - \boldsymbol{\mu}_{i}), \boldsymbol{\Omega}_{i}), \mathbf{\Gamma}_{i} = \boldsymbol{\Sigma}_{i}^{-1}\boldsymbol{B}_{i}, \boldsymbol{\Omega}_{i} = \boldsymbol{I}_{q} - \boldsymbol{B}_{i}^{\mathsf{T}}\boldsymbol{\Sigma}_{i}^{-1}\boldsymbol{B}_{i}, \\ \tilde{\pi}_{ij} = P(z_{ij} = 1 \mid \mathbf{v}_{j}, \mathbf{c}_{j}) = \frac{\pi_{i}f_{ij}(\mathbf{v}_{j} \mid \mathbf{c}_{j}, \boldsymbol{\theta}_{i})}{\sum_{t=1}^{g} \pi_{t}f_{tj}(\mathbf{v}_{j} \mid \mathbf{c}_{j}, \boldsymbol{\theta}_{t})} \end{aligned}$$

ML Estimation

> The EM algorithm

- Let $\boldsymbol{\Theta} = (\boldsymbol{\pi}_i, \{\boldsymbol{\beta}_i, \boldsymbol{B}_i, \boldsymbol{D}_i\}_{i=1}^g)$ be the entire unknown parameters in the model.
- **E-step:** $Q(\boldsymbol{\Theta} \mid \hat{\boldsymbol{\Theta}}^{(k)}) = E[\ell_c(\boldsymbol{\Theta} \mid \boldsymbol{y}, \boldsymbol{u}, \boldsymbol{z}) \mid \boldsymbol{v}_j, \boldsymbol{c}_j, \hat{\boldsymbol{\Theta}}^{(k)}]$
- **M-step:** Find $\hat{\Theta}^{(k+1)}$ by maximizing *Q*-function, yielding

 $\hat{\pi}_{i}^{(k+1)} = \frac{\sum_{j=1}^{n} \hat{z}_{ij}^{(k)}}{2}$ $\hat{\boldsymbol{\beta}}_{i}^{(k+1)} = \left(\sum_{i=1}^{n} \hat{z}_{ij}^{(k)} \boldsymbol{X}_{j}^{\top} \boldsymbol{X}_{j}\right)^{-1} \left(\sum_{i=1}^{n} \hat{z}_{ij}^{(k)} \boldsymbol{X}_{j}^{\top} (\hat{\boldsymbol{y}}_{ij}^{(k)} - \boldsymbol{B}_{i} \hat{\boldsymbol{u}}_{ij}^{(k)})\right),$ $\hat{\boldsymbol{B}}_{i}^{(k+1)} = \left(\sum_{i=1}^{n} \hat{z}_{ij}^{(k)} (\widehat{\boldsymbol{y}_{j}} \widehat{\boldsymbol{u}_{ij}}^{(k)} - \boldsymbol{X}_{j} \hat{\boldsymbol{\beta}}_{i}^{(k+1)} \hat{\boldsymbol{u}}_{ij}^{(k)^{\top}})\right) \left(\sum_{i=1}^{n} \hat{z}_{ij}^{(k)} \hat{\boldsymbol{u}}_{ij}^{2} \widehat{\boldsymbol{u}}_{ij}^{(k)}\right)^{-1},$ $\frac{\text{Diag}(\sum_{j=1}^{n} \hat{z}_{ij}^{(k)} \hat{\Psi}_{ij}^{(k+1)})}{\sum_{i=1}^{n} \hat{z}_{ij}^{(k)}}$

where

- > The AECM algorithm
- Partition the unknown parameters into two subjects: $\Theta_1 = \{\pi_i, \beta_i\}_{i=1}^g$ and $\Theta_2 = \{B_i, D_i\}_{i=1}^g$.
- The 1st cycle:
- **E-step:** $Q^{[1]}(\boldsymbol{\Theta}_1 \mid \hat{\boldsymbol{\Theta}}^{(k)}) = E[\ell_c^{[1]}(\boldsymbol{\Theta}_1 \mid \boldsymbol{y}, \boldsymbol{z}) \mid \boldsymbol{v}_j, \boldsymbol{c}_j, \hat{\boldsymbol{\Theta}}^{(k+1)}]$
- **CM-step:** Update $\hat{\Theta}_1^{(k+1)}$ by maximizing $Q^{[1]}$ function, yielding

 $\hat{\pi}_{i}^{(k+1)} = \frac{\sum_{j=1}^{n} \hat{z}_{ij}^{(k)}}{\hat{z}_{ij}}$ $\hat{\boldsymbol{\beta}}_{i}^{(k+1)} = \left(\sum_{i=1}^{n} \hat{z}_{ij}^{(k)} \boldsymbol{X}_{j}^{\top} \boldsymbol{X}_{j}\right)^{-1} \left(\sum_{i=1}^{n} \hat{z}_{ij}^{(k)} \boldsymbol{X}_{j}^{\top} \hat{\boldsymbol{Y}}_{ij}^{(k)}\right).$

- The 2nd cycle:
- **E-step:** $Q^{[2]}(\boldsymbol{\Theta}_2 \mid \hat{\boldsymbol{\Theta}}^{(k)}) = E[\ell_c^{[2]}(\boldsymbol{\Theta}_1 \mid \boldsymbol{y}, \boldsymbol{u}, \boldsymbol{z}) \mid \boldsymbol{v}_j, \boldsymbol{c}_j, \hat{\boldsymbol{\Theta}}^{(k+1)}]$
- **CM-step:** Update $\hat{\Theta}_2^{(k+1)}$ by maximizing $Q^{[2]}$ function, yielding
- The regression coefficients suggest that thirdgrade girls from urban zone in general have significantly better performance in speaking fluent Spanish than the other students.
- The two model selection criteria show that the data can be better described by fitting models with more than one component.



- The CPU time of EM and AECM algorithms are (1.8 min, 36 min) and (1.5 min, 35 min).
- All plots indicate the EM algorithm is worse than AECM in terms of the log-likelihood values and required iterations.
- The CPU time for the AECM algorithm is substantially less than those spent by EM, showing a dramatic advantage of using AECM for acceleration.

Conclusion & Future Work





Estimation of standard errors

The score vector and Hessian matrix

Let $\tilde{\Theta} = (\tilde{\theta}_1, \dots \tilde{\theta}_g)$, where $\theta_i = (\pi_i, \beta_i)$ represent the unknown parameters in the *i*th component. The individual score vector and Hessian matrix are defined as

$$s(\tilde{\Theta} \mid \boldsymbol{y}, \boldsymbol{z}) = \frac{\partial \ell_c(\tilde{\Theta} \mid \boldsymbol{y}, \boldsymbol{z})}{\partial \tilde{\Theta}} = (s_{\tilde{\theta}_1}^{\top}, \dots, s_{\tilde{\theta}_g}^{\top})^{\top} \text{ and } \boldsymbol{H}(\tilde{\Theta} \mid \boldsymbol{y}, \boldsymbol{z}) = \frac{\partial^2 \ell_c(\tilde{\Theta} \mid \boldsymbol{y}, \boldsymbol{z})}{\partial \tilde{\Theta} \partial \tilde{\Theta}^{\top}} = \text{Diag}\{\boldsymbol{H}_{\tilde{\theta}_1 \tilde{\theta}_1}, \dots, \boldsymbol{H}_{\tilde{\theta}_g \tilde{\theta}_g}\}, \text{ where}$$

$$s_{\tilde{\theta}_i} = \sum_{j=1}^n z_{ij} \boldsymbol{s}_{\tilde{\theta}_i}^{(j)} = \sum_{j=1}^n z_{ij} [\boldsymbol{s}_{\pi_i}^{(j)^{\top}}, \boldsymbol{s}_{\beta_i}^{(j)^{\top}}]^{\top}, \boldsymbol{s}_{\pi_i}^{(j)} = \frac{\partial \tilde{\ell}_{ij}}{\partial \pi_i} = \frac{1}{\pi_i}, \boldsymbol{s}_{\beta_i}^{(j)} = \frac{\partial \tilde{\ell}_{ij}}{\partial \beta_i} = \boldsymbol{X}_j^{\top} \boldsymbol{\Sigma}_i^{-1} (\boldsymbol{Y}_j - \boldsymbol{X}_j \beta_i),$$

$$H_{\tilde{\theta}_i \tilde{\theta}_i} = \sum_{j=1}^n z_{ij} \boldsymbol{H}_{\tilde{\theta}_i \tilde{\theta}_i}^{(j)} = \sum_{j=1}^n z_{ij} \begin{bmatrix} \boldsymbol{H}_{\pi_i \pi_i}^{(j)} & \boldsymbol{H}_{\pi_i \beta_i}^{(j)} \\ \boldsymbol{H}_{\beta_i \pi_i}^{(j)} & \boldsymbol{H}_{\beta_i \beta_i}^{(j)} \end{bmatrix}.$$
Furthermore, it can be verified that $\boldsymbol{H}_{\beta_i \pi_i}^{(j)} = \boldsymbol{H}_{\pi_i \beta_i}^{(j)^{\top}} = \boldsymbol{0}, \boldsymbol{H}_{\pi_i \pi_i}^{(j)} = \frac{\partial^2 \tilde{\ell}_{ij}}{\partial \pi_i^2} = -\pi_i^{-2},$
and $\boldsymbol{H}_{\beta_i \beta_i}^{(j)} = \frac{\partial^2 \tilde{\ell}_{ij}}{\partial \beta_i \partial \beta_i^{\top}} = -\boldsymbol{X}_j^{\top} \boldsymbol{\Sigma}_i^{-1} \boldsymbol{X}_j.$

 \succ Louis's formula(Louis, 1982): The Fisher information matrix of $\tilde{\Theta}$ can be approximated by $\boldsymbol{I}(\tilde{\boldsymbol{\Theta}} \mid \boldsymbol{y}) = -E\{\boldsymbol{H}(\tilde{\boldsymbol{\Theta}} \mid \boldsymbol{y}, \boldsymbol{z}) \mid \boldsymbol{v}, \boldsymbol{c}\} - \operatorname{cov}\{\boldsymbol{s}(\tilde{\boldsymbol{\Theta}} \mid \boldsymbol{y}, \boldsymbol{z}) \mid \boldsymbol{v}, \boldsymbol{c}\}\}$

- > We propose two feasible EM-type algorithms without resorting to direct evaluation of the intractable observed likelihood function for estimating the MFAC model with covariates.
- > Numerical results suggest the AECM algorithm with less amount of missing information has typically better convergence behavior than EM.
- > Future research would be interested in extending the current approach to a more flexible framework that can accommodate missing values and censored responses occurring simultaneously (Lin et al. 2018).

References

- Costa, D.R., Lachos, V.H., Bazan, J.L., Azevedo, C.L.N., (2014) Estimation methods for multivariate Tobit confirmatory factor analysis. Comput. Statist. Data Anal. 79, 248–260
- Hoffman, H., Johnson, R., (2015) Pseudo-likelihood estimation of multivariate normal parameters in the presence of leftcensored data. J. Agric. Biol. Environ. Stat. 20,156–171
- > Louis, T.A., (1982) Finding the observed information matrix when using the EM algorithm. J. Roy. Statist. Soc. B 44, 226– 233
- Example 2 Lin, T.I. (2009) Maximum likelihood estimation for multivariate skew normal mixture models. Journal of Multivariate Analysis 100, 257-265
- Ein, TI, Lachos, VH and Wang, WL (2018) Multivariate longitudinal data analysis with censored and intermittent missing responses. Statistics in Medicine 37, 2822-2835
- > Wang, W.L., Castro, L.M., Lachos, V.H., Lin, T.I. (2019) Model-based clustering of censored data via mixtures of factor analyzers, revised for Computational Statistics and Data Analysis