

Model-Based Clustering Analysis for Multivariate Longitudinal Data via Mixtures of Matrix-Variate t-distributions

Farzane Ahmadi^{1*}, Elham Faghihzadeh²

¹ Department of Biostatistics and Epidemiology, Faculty of Medicine, Zanjan University of Medical Sciences, Zanjan, Iran
ahmadi.farzane@zums.ac.ir; faghihzadeh.elham@gmail.com

² Independent researcher, Tehran, Iran

Abstract: The finite mixture model is considered as an appropriate instrument for data clustering. Different parsimonious multivariate mixture distributions are introduced for skewed and/or heavy-tailed longitudinal data. The eigenvalue or modified Cholesky decomposition of covariance matrices develops the families of parsimonious mixture models. Thus, the finite mixture of matrix-variate t-distributions for clustering a three-way dataset with heavy-tailed or outlier observations (e.g., multivariate longitudinal data) is more appropriate compared to matrix-variate normal distributions. Accordingly, the present study considered a parsimonious family of the finite mixture of matrix-variate t-distributions using the eigenvalue and modified Cholesky decomposition for within and between covariance matrices, respectively. Finally, parameter estimates were calculated using the expectation-maximization algorithm, and simulations studies and real data analyses were conducted to confirm the obtained results.

Keywords: Eigenvalue Decomposition; Finite Mixture; Matrix-Variate t-Distribution; Modified Cholesky Decomposition; Multivariate Longitudinal Data; Parsimonious Covariance Structures

1. Introduction

Finite mixture models in the statistical data analysis mainly contribute to modelling a heterogeneous population and providing an easy and model-based method for clustering and classification structure [1], [2]. Different studies have evaluated various finite mixtures of distributions focusing on multivariate (two-way data) distributions. For instance, such studies have proposed different finite mixtures of multivariate distributions, including multivariate normal distribution [3], multivariate t-distribution [4], multivariate skew-normal distribution [5], multivariate skew-t-distribution [6], multivariate normal inverse Gaussian [7], multivariate generalized hyperbolic distribution [8], and multivariate power exponential distribution [9] over the last two decades.

Three-way data including multivariate longitudinal, spatial multivariate, and spatio-temporal data may be available in a range of scientific domains [10]. Despite the important role of matrix-variate distributions in three-way data analysis, a small body of research exists in this respect. For example, *Viroli* introduced the finite mixtures of matrix-variate normal distributions (MVNDs) for classifying the three-way data [11]. In addition, *Anderlucci* and *Viroli* [12] considered the finite mixture of MVNDs for multivariate longitudinal data. In another study, *Doğru*, *Bulut* and *Arslan* [13] proposed a finite mixture of matrix-variate t-distributions (MVTDs). Further, *Gallaugh* and *McNicholas* [14]–[16] applied four skewed matrix-variate distributions of matrix-variate skew-t, generalized hyperbolic, variance-gamma, and normal inverse Gaussian distributions in the finite mixture of these distributions. In the two- or three-way data, where there are some departures from normality in datasets, using normal distributions affects the estimation of some parameters [17]. The presence of outlier or heavy-tailed data is considered as one of the common departures from normality and in such case, the mixture of t-distributions is an appropriate alternative to the mixture of normal distributions [13].

On the other hand, without any constraints on mixture parameters, the number of estimated parameters increases dramatically by an increase in components. Therefore, some constraints should be put on model parameters in order to achieve more parsimonious models. Considering a large number of mixture parameters in the covariance matrix component, more attention should be drawn on covariance structure decomposition. Further literature contains parsimonious covariance matrices in the mixture of multivariate distributions [18]–[24]. Some studies have investigated the parsimonious feature only in the finite mixture of MVNDs for three-way data [11], [12]. However, to the best of our knowledge, no research has applied the parsimonious MVTD mixture model. Therefore, the present study focused on the parsimonious mixture of MVTDs for

clustering multivariate longitudinal data with outliers or heavy-tails. The remaining sections of the present study are organized as follows. Section 2 reviews the finite mixture of MVTDs, the decomposition of covariance matrices, and the details of the estimates of MVTD parameters. Furthermore, Section 3 discusses the simulation studies and real examples in order to demonstrate the performance of models.

2. Method

2.1 Background

2.1.1 Finite mixture of MVTDs

A $T \times p$ dimensional random matrix \mathbf{X} is assumed to arise from a parametric finite mixture if it is possible to write $p(\mathbf{X}|\boldsymbol{\vartheta}) = \sum_{i=1}^k \pi_i p_i(\mathbf{X}|\theta_i)$ for all $\mathbf{X} \in \mathcal{X}$, where $\boldsymbol{\vartheta} = (\pi_1, \dots, \pi_k, \theta_1, \dots, \theta_k)$ is the vector of parameters, and π_i and k are the mixing proportion and the number of mixture components, respectively, so that $\sum_{i=1}^k \pi_i = 1$ and $\pi_i \in [0,1]$. Additionally, $p_i(\mathbf{X}|\theta_i)$ is referred to as the density of the component. In the mixture of MVTDs, component density with a $T \times p$ mean matrix \mathbf{M}_i , two covariance matrices $\boldsymbol{\Phi}_i$ and $\boldsymbol{\Omega}_i$ with dimensions $T \times T$ and $p \times p$, and degrees of freedom v_i is as follows [25]:

$$Mt^{(T \times p)}(\mathbf{X}|\mathbf{M}_i, \boldsymbol{\Phi}_i, \boldsymbol{\Omega}_i, v_i) = \frac{\Gamma\left(\frac{Tp + v_i}{2}\right)}{(\pi v_i)^{\frac{Tp}{2}} \Gamma\left(\frac{v_i}{2}\right)^{\frac{p}{2}} |\boldsymbol{\Phi}_i|^{\frac{T}{2}} |\boldsymbol{\Omega}_i|^{\frac{T}{2}}} \left(1 + \frac{\text{tr}\{(\mathbf{X} - \mathbf{M}_i)' \boldsymbol{\Phi}_i^{-1} (\mathbf{X} - \mathbf{M}_i) \boldsymbol{\Omega}_i^{-1}\}}{v_i}\right)^{-\frac{Tp+v_i}{2}} \quad (1)$$

where T and p indicate the number of measurement times and the number of response variables, respectively. In addition, $\boldsymbol{\Phi}_i$ and $\boldsymbol{\Omega}_i$ are commonly referred to as *between* and *within* covariance matrices, respectively. In the present study, the *upper case boldface* was used for the matrices. The MVTDs arise as a particular case of a normal-variance mixture distribution. In this formulation, the random matrix \mathbf{X} is defined as $\mathbf{X} = \mathbf{M} + W^{\frac{1}{2}}\mathbf{V}$, where the matrix random \mathbf{V} has the MVND with the mean matrix $\mathbf{0}$ and covariance matrices $\boldsymbol{\Phi}$ and $\boldsymbol{\Omega}$, $\mathbf{V} \sim \phi^{(T \times p)}(\mathbf{X}|\mathbf{M}, \boldsymbol{\Phi}, \boldsymbol{\Omega})$, and the latent random variable W follows an inverse gamma distribution with parameters $\left(\frac{v}{2}, \frac{v}{2}\right)$ [13]. In addition, the estimates of $\boldsymbol{\Omega}_i$ and $\boldsymbol{\Phi}_i$ are not unique. For each positive and nonzero constant a , we have $\boldsymbol{\Omega}_i \otimes \boldsymbol{\Phi}_i = a\boldsymbol{\Omega}_i \otimes \left(\frac{1}{a}\boldsymbol{\Phi}_i\right)$. The constraint $\text{tr}(\boldsymbol{\Omega}_i) = p$ or $\text{tr}(\boldsymbol{\Phi}_i) = T$ can be used to obtain an identifiable solution for $\boldsymbol{\Omega}_i$ and $\boldsymbol{\Phi}_i$ [12], [14].

2.1.2 The decomposition of covariance matrices

Restrictions on mixture parameters are typically constructed by constraining covariance matrices. To achieve parsimonious models, eigenvalue and the modified Cholesky decompositions were used for the *between* and *within* covariance matrices, respectively.

The eigenvalue decomposition

The eigenvalue decomposition was used in multivariate normal mixtures and the other multivariate mixture distributions such as t-mixture distributions [22], along with skew-normal and skew-t mixture distributions [24] for clustering, classification, and discriminant analysis. On the other hand *Viroli* [11] and *Sarkar et al.* [27] applied the eigenvalue decomposition in the mixture of MVNDs. This parameterization includes the expression *within* component-covariance matrix ($\boldsymbol{\Omega}_i$) in terms of its eigenvalue decomposition as $\boldsymbol{\Omega}_i = \lambda_i \mathbf{D}_i \mathbf{A}_i \mathbf{D}_i'$, where \mathbf{D}_i denotes the matrix of eigenvectors. Furthermore, \mathbf{A}_i is a diagonal matrix whose elements are proportional to the eigenvalues of $\boldsymbol{\Omega}_i$ and λ_i represents the associated proportionality constant. Different sub-models can be defined by considering homoscedastic or varying quantities across mixture components [11], [20].

Modified Cholesky decomposition

The *between* component-covariance matrix (Φ_i) of the multivariate longitudinal data can be decomposed by the modified Cholesky decomposition. *Anderlucci* and *Viroli* [12] employed the above-mentioned decomposition along with the eigenvalue decomposition for the *between* and *within* covariance structures in the mixture of MVNDs, respectively. The modified Cholesky decomposition was expressed as $\Phi_i^{-1} = \mathbf{U}_i' \mathbf{T}_i^{-1} \mathbf{U}_i$ where \mathbf{U}_i is a unique lower triangular matrix with diagonal elements 1 and \mathbf{T}_i denotes a unique diagonal matrix with strictly positive diagonal entries representing innovation variances. The lower diagonal elements in \mathbf{U}_i equal the negative coefficients resulted from the regression of X_t on $X_{t-1}, X_{t-2}, \dots, X_1$ e.g. $\hat{X}_t = M_t + \sum_{s=1}^{t-1} \phi_{r,s}^{(i)} (X_t - M_t)$ [28]. On the other hand, different orders (m) can be considered in matrix \mathbf{U}_i , where m can range from 0 to T-1. The lower orders provide more parsimonious models so that m=0 and m=1 equal the independency of different times and the dependency of X_t on a previous time (X_{t-1}), and the like. Accordingly, the modified Cholesky decomposition for the temporal covariance matrix equals the generalized autoregressive with process order m, GAR(m). Thus, the r^{th} row elements of matrix \mathbf{U}_i which should be estimated can be written as $(-\phi_{r,r-m}^{(i)}, -\phi_{r,r-m+1}^{(i)}, \dots, -\phi_{r,r-1}^{(i)})'$; $r = 2, \dots, T$, $m = 0, 1, \dots, T-1$. Additionally, matrix \mathbf{T}_i can be defended as $\mathbf{T}_i = d_i I_T$ (Isotropic) for a more parsimonious model. In addition, different sub-models can be defined by considering homoscedastic or varying quantities (i.e., \mathbf{U}_i and \mathbf{T}_i) across mixture components [12].

2.2 Estimation of parameters

To find the maximum likelihood estimators for mixture parameters, the present study used an expectation-maximization (EM) algorithm [29] which was proposed by *Dođru*, *Bulut* and *Arslan* [13] for the mixture of MVTDs. Assume that $\mathbf{X}_1, \dots, \mathbf{X}_n$, where n is the number of observations, be a random sample of matrices from the mixture of MVTDs, and Z_{ij} denotes the component membership of observation j . Further, $Z_{ij} = 1$ if the j^{th} observation is from component i , otherwise, $Z_{ij} = 0$, where $j = 1, \dots, n$ and $i = 1, \dots, k$. MVTDs are expressed as follows:

$$\mathbf{X}_j | W_j, Z_{ij} = 1 \sim \phi^{(T \times p)}(\mathbf{M}_i, W_j^{-1} \Phi_i, \Omega_i), \quad W_j | Z_{ij} = 1 \sim \text{Gamma}\left(\frac{v_i}{2}, \frac{v_i}{2}\right). \quad (2)$$

Based on the hierarchical representation of the MVTDs, the complete data log-likelihood $\ell_c(\vartheta)$ can be written as follows:

$$\begin{aligned} \ell_c(\vartheta) = & \sum_{j=1}^n \sum_{i=1}^k Z_{ij} \left[-\frac{Tp}{2} \log 2\pi - \frac{p}{2} \log |u_j^{-1} \Phi_i| - \frac{T}{2} \log |\Omega_i| - \frac{W_j}{2} \text{tr} \left\{ \Omega_i^{-1} (\mathbf{X}_j - \mathbf{M}_i) \Phi_i^{-1} (\mathbf{X}_j - \mathbf{M}_i)' \right\} + \frac{v_i}{2} \log \left(\frac{v_i}{2} \right) \right. \\ & \left. - \log \Gamma \left(\frac{v_i}{2} \right) - \frac{v_i}{2} W_j + \left(\frac{v_i}{2} - 1 \right) \log(W_j) \right] + \sum_{j=1}^n \sum_{i=1}^k Z_{ij} \log \pi_i. \end{aligned} \quad (3)$$

An EM algorithm is as follows:

- I. **Initialization:** Initialize parameters π_i , \mathbf{M}_i , Φ_i , Ω_i , and v_i , setting $t = 0$.
- II. **E-step:** Update $E(Z_{ij} | \mathbf{X}_j, \vartheta)$, $E(W_j | \mathbf{X}_j, Z_{ij} = 1; \vartheta)$, and $E(\log W_j | \mathbf{X}_j, Z_{ij} = 1; \vartheta)$, where

$$\begin{aligned} E(Z_{ij} | \mathbf{X}_j, \vartheta^{(t)}) = P(Z_{ij} = 1 | \mathbf{X}_j, \vartheta^{(t)}) &= \frac{\pi_i^{(t)} M t^{T \times p}(\mathbf{X}_j; \mathbf{M}_i^{(t)}, \Phi_i^{(t)}, \Omega_i^{(t)}, v_i^{(t)})}{\sum_{i=1}^k \pi_i^{(t)} M t^{T \times p}(\mathbf{X}_j; \mathbf{M}_i^{(t)}, \Phi_i^{(t)}, \Omega_i^{(t)}, v_i^{(t)})} = \tau_{ij}^{(t)}, \\ E(W_j | \mathbf{X}_j, Z_{ij} = 1; \vartheta^{(t)}) &= \frac{Tp + \hat{v}_i^{(t)}}{\text{tr} \left\{ \Omega_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t)})' \Phi_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t)}) \right\} + v_i^{(t)}} = W_{1ij}^{(t)}, \\ (\log W_j | \mathbf{X}_j, Z_{ij} = 1; \vartheta^{(t)}) &= DG \left(\frac{Tp + v_i^{(t)}}{2} \right) + \log \left(\frac{\text{tr} \left\{ \Omega_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t)})' \Phi_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t)}) \right\} + v_i^{(t)}}{2} \right) = W_{2ij}^{(t)}, \end{aligned} \quad (4)$$

where $DG(t) = \frac{d}{dt} \log \Gamma(T)$ represents the digamma function. Furthermore, $W_{1ij}^{(t)}$ is calculated based on $W_j | X_j, Z_{ij} = 1$ distribution, which has a gamma distribution with parameters $\frac{T p + \hat{v}_i^{(t)}}{2}$ and $\frac{\text{tr}\{\boldsymbol{\Omega}_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t)})' \boldsymbol{\Phi}_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t)})\} + v_i^{(t)}}{2}$, and $W_{2ij}^{(t)}$ is achieved using the moment-generating function of $W_j | X_j, Z_{ij} = 1$.

III. M-step: Update $\pi_i, \mathbf{M}_i, \boldsymbol{\Omega}_i, \boldsymbol{\Phi}_i$, and v_i . The order of parameter estimation is as follows (1): π_i and \mathbf{M}_i ; (2) $\boldsymbol{\Omega}_i$; (3) $\boldsymbol{\Phi}_i$; (4) v_i

$$1. \text{ Update } \pi_i \text{ and } \mathbf{M}_i: \quad \pi_i^{(t+1)} = \frac{\sum_{j=1}^n \tau_{ij}^{(t)}}{n}, \quad \mathbf{M}_i^{(t+1)} = \frac{\sum_{j=1}^n \tau_{ij}^{(t)} W_{1ij}^{(t)} \mathbf{X}_j}{\sum_{j=1}^n \tau_{ij}^{(t)} W_{1ij}^{(t)}}, \quad (5)$$

2. Update $\boldsymbol{\Omega}_i$

Assuming that $\mathbf{B}_i = \sum_{j=1}^n \tau_{ij}^{(t)} W_{1ij}^{(t)} (\mathbf{X}_j - \mathbf{M}_i^{(t+1)})' \boldsymbol{\Phi}_i^{(t)-1} (\mathbf{X}_j - \mathbf{M}_i^{(t+1)})$, the $\ell_c(\vartheta)$ is proportional to $-\frac{T}{2} \sum_{i=1}^k n_i \log |\boldsymbol{\Omega}_i| - \frac{1}{2} \sum_{i=1}^k \text{tr}\{\boldsymbol{\Omega}_i^{-1} \mathbf{B}_i\}$ with $n_i = \sum_{j=1}^n \tau_{ij}^{(t)}$. The estimates of parameters for the eight sub-models are provided below.

- *Sub-model VVV:* The maximization of $-\frac{T}{2} \sum_{i=1}^k n_i \log |\boldsymbol{\Omega}_i| - \frac{1}{2} \sum_{i=1}^k \text{tr}\{\boldsymbol{\Omega}_i^{-1} \mathbf{B}_i\}$ with respect to $\boldsymbol{\Omega}_i$ leads to $\boldsymbol{\Omega}_i^{(t+1)} = \frac{\mathbf{B}_i}{n_i T}$;
- *Sub-model EEE:* The maximization of $-\frac{Tn}{2} \log |\boldsymbol{\Omega}| \sum_{i=1}^k n_i - \frac{1}{2} \text{tr}\{\boldsymbol{\Omega}^{-1} \sum_{i=1}^k \mathbf{B}_i\}$, where $n = \sum_{i=1}^k \sum_{j=1}^n \tau_{ij}^{(t)}$, with respect to $\boldsymbol{\Omega}_i = \boldsymbol{\Omega}$ leads to $\boldsymbol{\Omega}^{(t+1)} = \frac{\sum_{i=1}^k \mathbf{B}_i}{nT}$;
- *Sub-model VVI:* The maximization of $-\frac{Tp}{2} \sum_{i=1}^k n_i \log \lambda_i - \frac{1}{2} \sum_{i=1}^k \frac{1}{\lambda_i} \text{tr}\{\mathbf{A}_i^{-1} \mathbf{B}_i\}$ with respect to $\boldsymbol{\Omega}_i = \lambda_i \mathbf{D} \mathbf{A}_i \mathbf{D}'$ leads to $\lambda_i^{(t+1)} = \frac{|\text{diag}(\mathbf{B}_i)|^{\frac{1}{p}}}{T n_i}$ and $\mathbf{A}_i^{(t+1)} = \frac{\text{diag}(\mathbf{B}_i)}{|\text{diag}(\mathbf{B}_i)|^{\frac{1}{p}}}$;
- *Sub-model EEI:* The maximization of $-\frac{pT}{2} \sum_{i=1}^k n_i \log |\boldsymbol{\Omega}_i| - \frac{1}{2} \sum_{i=1}^k \text{tr}\{\boldsymbol{\Omega}_i^{-1} \mathbf{B}_i\}$ with respect to $\boldsymbol{\Omega}_i$ leads to $\lambda^{(t+1)} = \frac{|\text{diag}(\sum_{i=1}^k \mathbf{B}_i)|^{\frac{1}{p}}}{Tn}$ and $\mathbf{A}^{(t+1)} = \frac{\text{diag}(\sum_{i=1}^k \mathbf{B}_i)}{|\text{diag}(\sum_{i=1}^k \mathbf{B}_i)|^{\frac{1}{p}}}$;
- *Sub-model VII:* The maximization of $-\frac{pT}{2} \sum_{i=1}^k n_i \log \lambda_i - \frac{1}{2} \sum_{i=1}^k \text{tr}\left\{\frac{\mathbf{B}_i}{\lambda_i}\right\}$ with respect to $\boldsymbol{\Omega}_i = \lambda_i \mathbf{I}$ leads to $\lambda_i^{(t+1)} = \frac{\text{tr}\{\mathbf{B}_i\}}{T p n_i}$;
- *Sub-model EII:* The maximization of $-\frac{Tpn}{2} \log \lambda - \frac{1}{2\lambda} \text{tr}\{\sum_{i=1}^k \mathbf{B}_i\}$ with respect to $\boldsymbol{\Omega} = \lambda \mathbf{I}$ leads to $\lambda^{(t+1)} = \frac{\text{tr}\{\sum_{i=1}^k \mathbf{B}_i\}}{Tpn}$;
- *Sub-model EEV:* The maximization of $-\frac{Tpn}{2} \log \lambda - \frac{1}{2\lambda} \sum_{i=1}^k \text{tr}\{\mathbf{D}_i \mathbf{A}^{-1} \mathbf{D}_i' \mathbf{B}_i\}$ with respect to $\boldsymbol{\Omega}_i = \lambda \mathbf{D}_i \mathbf{A} \mathbf{D}_i'$ leads to $\lambda^{(t+1)} = \frac{|\sum_{i=1}^k \mathbf{C}_i|^{\frac{1}{p}}}{nT}$, $\mathbf{A}^{(t+1)} = \frac{\sum_{i=1}^k \mathbf{C}_i}{|\sum_{i=1}^k \mathbf{C}_i|^{\frac{1}{p}}}$, $\mathbf{D}_i^{(t+1)} = \mathbf{L}_i$, where for $i = 1, \dots, k$ \mathbf{C}_i , and \mathbf{L}_i are derived from the eigenvalue decomposition of the symmetric positive definite matrix $\mathbf{B}_i = \mathbf{L}_i \mathbf{C}_i \mathbf{L}_i'$ with the eigenvalues in the diagonal matrix \mathbf{C}_i in descending order.
- *Sub-model III:* This situation equals the independence of the responses thus no parameters are available.

3. Update $\boldsymbol{\Phi}_i$

Considering that $\mathbf{S}^{(i)} = \sum_{j=1}^n \tau_{ij}^{(t)} W_{1ij}^{(t)} (\mathbf{X}_j - \hat{\mathbf{M}}_i^{(t+1)}) \boldsymbol{\Omega}_i^{(t+1)-1} (\mathbf{X}_j - \hat{\mathbf{M}}_i^{(t+1)})'$, $\ell_c(\vartheta)$ is proportional to $-\frac{p}{2} \sum_{i=1}^k n_i \log |\mathbf{D}_i| - \frac{1}{2} \text{tr}\{\sum_{i=1}^k (\mathbf{U}_i' \mathbf{T}_i^{-1} \mathbf{U}_i) \mathbf{S}^{(i)}\}$. The estimates of parameters for the four sub-models are presented as follows:

- *Sub-model GAR(m):* The maximization of $-\frac{p}{2} \sum_{i=1}^k n_i \log |\mathbf{D}_i| - \frac{1}{2} \text{tr}\{\sum_{i=1}^k (\mathbf{U}_i' \mathbf{T}_i^{-1} \mathbf{U}_i) \mathbf{S}^{(i)}\}$ with respect to $\boldsymbol{\Phi}_i$ leads to the r^{th} row estimation of matrix \mathbf{U}_i as

$$\begin{pmatrix} \phi_{r,r-m}^{(i)} \\ \phi_{r,r-m+1}^{(i)} \\ \dots \\ \phi_{r,r-1}^{(i)} \end{pmatrix}^{(t+1)} = \begin{pmatrix} S_{r-m,r-m}^{(i)} & S_{r-m+1,r-m}^{(i)} & \dots & S_{r-1,r-m}^{(i)} \\ S_{r-m,r-m+1}^{(i)} & S_{r-m+1,r-m+1}^{(i)} & \dots & S_{r-1,r-m+1}^{(i)} \\ \dots & \dots & \ddots & \vdots \\ S_{r-m,r-1}^{(i)} & S_{r-m+1,r-1}^{(i)} & \dots & S_{r-1,r-1}^{(i)} \end{pmatrix}^{-1} \begin{pmatrix} S_{r,r-m}^{(i)} \\ S_{r,r-m+1}^{(i)} \\ \dots \\ S_{r,r-1}^{(i)} \end{pmatrix}, \quad (6)$$

and matrix $\mathbf{T}_i^{(t+1)} = \frac{1}{p} \text{diag}(\mathbf{U}_i^{(t+1)} \mathbf{S}^{(i)} \mathbf{U}_i^{(t+1)'})$, where $r = 2, \dots, T$, $m = 0, \dots, T-1$, and $S_{l,t}^{(i)}$ is the l^{th} -row and t^{th} -column element of matrix $\mathbf{S}^{(i)}$.

- *Sub-model GARI(m)*: The maximization of $-\frac{Tp}{2} \sum_{i=1}^k n_i \log |d_i| - \frac{1}{2} \text{tr} \left\{ \sum_{i=1}^k \frac{1}{d_i} \mathbf{U}'_i \mathbf{U}_i \mathbf{S}^{(i)} \right\}$ with respect to $\Phi_i = \frac{1}{d_i} \mathbf{U}'_i \mathbf{U}_i$ leads to the same estimate of \mathbf{U}_i as *sub-model GAR(m)* and estimate $d_i^{(t+1)} = \frac{\text{tr} \{ \mathbf{U}_i^{(t+1)} \mathbf{S}^{(i)} \mathbf{U}_i^{(t+1)'} \}}{n_i p T}$.
- *Sub-model EGAR(m)*: The maximization of $-\frac{np}{2} \log |\mathbf{D}| - \frac{1}{2} \text{tr} \{ \mathbf{U}' \mathbf{T}^{-1} \mathbf{U} (\sum_{i=1}^k \mathbf{S}^{(i)}) \}$ with respect to $\Phi_i = \Phi$ leads to the same estimate of \mathbf{U} as *sub-model GAR(m)* by replacing $\sum_{i=1}^k \mathbf{S}^{(i)}$ instead of $\mathbf{S}^{(i)}$ and estimate $\mathbf{T}^{(t+1)} = \frac{1}{np} \text{diag}(\mathbf{U}^{(t+1)} \{ \sum_{i=1}^k \mathbf{S}^{(i)} \} \mathbf{U}^{(t+1)'})$.
- *Sub-model EGARI(m)*: The maximization of $-\frac{npT}{2} \log |d| - \frac{1}{2d} \text{tr} \{ \mathbf{U}' \mathbf{U} (\sum_{i=1}^k \mathbf{S}^{(i)}) \}$ with respect to $\Phi_i = \Phi$ leads to the same estimate of \mathbf{U} as *sub-model EGAR(m)* and estimate $d^{(t+1)} = \frac{\text{tr} \{ \mathbf{U}^{(t+1)'} \mathbf{U}^{(t+1)} \sum_{i=1}^k \mathbf{S}^{(i)} \}}{pnT}$.

4. Update v_i

For the degree of freedom, two situations were considered, including equal and unequal v_i across mixture components (constrained and unconstrained v_i , respectively). Given $\tau_{ij}^{(t+1)}$, $\pi_i^{(t+1)}$, $\mathbf{M}_i^{(t+1)}$, $\mathbf{\Omega}_i^{(t+1)}$, and $\Phi_i^{(t+1)}$, the estimations of v_i are calculated by finding the root of equations (8) and (9) in constrained and unconstrained situations, respectively.

$$1 + \log \left(\frac{v}{2} \right) - DG \left(\Gamma \left(\frac{v}{2} \right) \right) + \frac{1}{n} \sum \sum \tau_{ij}^{(t+1)} (W_{2ij}^{(t+1)} - W_{1ij}^{(t+1)}) = 0, \quad (7)$$

$$1 + \log \left(\frac{v_i}{2} \right) - DG \left(\Gamma \left(\frac{v_i}{2} \right) \right) + \frac{1}{n_i} \sum \tau_{ij}^{(t+1)} (W_{2ij}^{(t+1)} - W_{1ij}^{(t+1)}) = 0. \quad (8)$$

IV. Check the convergence criterion: If not satisfied, set $t = t + 1$ and go to step **II** of the EM algorithm iteration.

2.3 Model selection and convergence criterion

It is possible to define a large family ($64 \times T$) of possible mixture models by allowing different sub-models for covariance matrices $\mathbf{\Omega}_i$ and Φ_i with different orders for matrix \mathbf{T}_i , $m = 0, 1, \dots, T-1$, and constrained/unconstrained for v_i . The model can be selected according to the Bayesian information criterion (BIC) as $BIC = 2 l(x, \hat{\theta}) - h \log n$, where $l(x, \hat{\theta})$ and $\hat{\theta}$ indicate the maximized log-likelihood and the maximum likelihood estimate of θ , respectively. Additionally, h and n are the number of free parameters in the model and the number of observations, respectively [30]. Other criteria are employed in addition to BIC to estimate the number of mixture components, such as Integrated Completed Likelihood (ICL), which is computed as $ICL \approx BIC - 2 \sum_{j=1}^n \sum_{i=1}^k \text{MAP}^1(\tau_{ij}^{(t)}) \log \tau_{ij}^{(t)}$, where $\text{MAP}(\tau_{ij}^{(t)}) = 1$ if the $\max_{i=1, \dots, k} \{ \tau_{ij}^{(t)} \} = i$, otherwise, $\text{MAP}(\tau_{ij}^{(t)}) = 0$, $j = 1, \dots, n$ and $i = 1, \dots, k$ [31].

¹ Maximum A Posteriori probability (MAP)

In general, 20 random multistate points were considered given that the starting values of the EM algorithm could affect the estimated parameters. If the convergence criterion $|l(x, \hat{\theta}^{(t+1)}) - l(x, \hat{\theta}^{(t)})| < 1.0e - 5$ is met, the EM algorithm is stopped, and the range of values for $\hat{\nu}_i$ is limited to between 2 and 200 [22]. These models have been written in R software.

3 Results: Simulation studies and real data

3.1 Simulation 1

The first simulation study was conducted to evaluate the ability of the algorithm for recognizing the temporal structure. The features of simulation study 1 were: a number k of mixture components equal to 3, a k -vector of the degrees of freedom equal to 5, 5, 5, and a 4×4 *within* covariance matrix Ω_i with a structure equals to VVV. In addition, other features included a 6×6 temporal covariance matrix Φ_i with a structure equals to GAR(1) and GAR(3), and a sample size n equals 100, 200, 500, and 1000. For each setting, 100 datasets were generated from the mixture of the MVTDs based on the defined *within* and temporal covariance matrices. Then, the mixture of MVTDs and MVNDs was run for five different models according to different orders for Φ_i : GAR(1), GAR(2), GAR(3), GAR(4), and GAR(5). The best model was chosen according to the BIC and ICL. The percentage (number) of correct model selection with MVTD is equal to 100 in all cases, and it ranges from 97 to 100 for MVNDs. In a situation with a true model GAR(3), this percentage varied from 99 to 100 and 93 to 99 for MVTD and MVND, respectively.

3.2 Simulation 2

It was performed to evaluate the ability of MVTDs to recover the MVNDs in multivariate longitudinal data. To this end, datasets were generated from two-component ($k=2$), matrix-variate mixture models. The MVND and MVTD were the first and second components, respectively, and the same covariance structures with different parameter values were used accordingly. Other features (i.e., Ω_i , Φ_i , and sample size) in this simulation are similar to the first simulation study. For each setting, 100 datasets were generated from mixture distributions based on the defined *within* and temporal covariance matrices. Further, the mixture of MVTDs and MVNDs was run for five different models of GAR(1), GAR(2), GAR(3), GAR(4), and GAR(5). Table 5 presents the average values of the degree of freedom (standard deviation) of a model with GAR(1) and GAR(3) structures for Φ_i . Given k ($=2$) and Ω_i (VVV), the estimated degrees of freedom demonstrated that the first component was normal. Furthermore, the degrees of freedom estimates were computed to be close to true values in MVTD mixture models.

Additionally, the misclassification error rate (MISC) and the measure of accuracy (γ) for mean and covariance matrices were computed for each dataset and model in order to compare the two models in parameter estimates. Therefore, the accuracy measures of \mathbf{M} , Ω (=VVV), \mathbf{U} , and \mathbf{T} (=GAR) were calculated by the following expressions [32]:

$$\gamma_M = \frac{\sum_{i=1}^k \|\hat{\mathbf{M}}_i - \mathbf{M}_i\|}{kTp}, \quad \gamma_\Omega = \frac{\sum_{i=1}^k \|\hat{\Omega}_i - \Omega_i\|}{\frac{kp(p+1)}{2}}, \quad \gamma_U = \frac{\sum_{i=1}^k \|\hat{\mathbf{U}}_i - \mathbf{U}_i\|}{k\phi}, \quad \gamma_T = \frac{\sum_{i=1}^k \|\hat{\mathbf{T}}_i - \mathbf{T}_i\|}{kT} \quad (9)$$

where the lower accuracy measure (γ) implies higher accuracy for parameters.

Considering k ($=2$) and Ω_i (VVV), γ_M , γ_Ω , and γ_T were not sensitive to the misspecification of the order of the temporal covariance ($m=1, 2, \dots, 5$), and these values were nearly identical in MVTD and MVND mixture models. However, the values of γ_T relied on the misspecification of the temporal covariance order. In the two models, γ_T of the lower orders ($m = 1, 2$) was larger compared to the higher orders ($m = 3, 4, 5$). It should be noted that MVND mixture models tend to overestimate γ_T compared to MVTD mixture models. Eventually, the accuracy measures in both models decreased by an increase in the sample size (Table 1). The mean compute time for fitting the mixture of MVTDs vs as MVNDs with the true model GAR(3) was 6.66 vs. 0.37 for $n=100$, 10.17 vs. 0.62 for $n=200$, 23.44 vs. 1.30 for $n=500$, and 46.86 vs 2.59 for $n=1000$.

Table 1: Mean (S.D) of MISC and accuracy measures with GAR(3) structure for Φ_i from simulation 2

		Φ_i									
		GAR(1)		GAR(2)		GAR(3)		GAR(4)		GAR(5)	
n		MVTD	MVND	MVTD	MVND	MVTD	MVND	MVTD	MVND	MVTD	MVND
200	MISC	0	0.0001(0.001)	0	0.0001(0.001)	0	0.0001(0.001)	0	0.0001(0.001)	0	0.0001(0.001)
	γ_M	0.09 (0.01)	0.10 (0.01)	0.09 (0.01)	0.10 (0.01)	0.09 (0.01)	0.10 (0.01)	0.09 (0.01)	0.10 (0.01)	0.09 (0.01)	0.10 (0.01)
	γ_Ω	0.23 (0.002)	0.23 (0.003)	0.23 (0.001)	0.23 (0.002)	0.23 (0.001)	0.23 (0.001)	0.23 (0.001)	0.23 (0.001)	0.23 (0.001)	0.23 (0.001)
	γ_T	0.49 (0.02)	0.57 (0.06)	0.44 (0.02)	0.52 (0.13)	0.40 (0.02)	0.48 (0.14)	0.40 (0.02)	0.48 (0.14)	0.40 (0.02)	0.48 (0.14)
	γ_U	0.33 (0.002)	0.33 (0.002)	0.20 (0.01)	0.20 (0.02)	0.10 (0.02)	0.11 (0.03)	0.15 (0.04)	0.17 (0.05)	0.18 (0.04)	0.21 (0.06)
	MISC	0	0	0	0	0	0	0	0	0	0
1000	γ_M	0.06 (0.01)	0.07 (0.01)	0.06 (0.01)	0.07 (0.01)	0.06 (0.01)	0.07 (0.01)	0.06 (0.01)	0.07 (0.01)	0.06 (0.01)	0.07 (0.01)
	γ_Ω	0.15(0.0002)	0.15 (0.0002)	0.15(0.0002)	0.15 (0.0002)	0.15(0.0002)	0.15 (0.0002)	0.15(0.0002)	0.15 (0.0002)	0.15(0.0002)	0.15 (0.0002)
	γ_T	0.31 (0.01)	0.42 (0.03)	0.27 (0.01)	0.40 (0.02)	0.27 (0.01)	0.37 (0.02)	0.27 (0.01)	0.37 (0.02)	0.27 (0.01)	0.37 (0.02)
	γ_U	0.28 (0.001)	0.28 (0.001)	0.17 (0.001)	0.20 (0.001)	0.05 (0.001)	0.06 (0.01)	0.05 (0.001)	0.06 (0.01)	0.05 (0.001)	0.06 (0.01)
	MISC	0	0	0	0	0	0	0	0	0	0

3.3 Simulation 3

The ability of MVTD and the MVND mixture models was evaluated regarding recognizing the true number of mixture components when the data were generated from MVTD mixture models. Then, the impact of the misspecification of the temporal matrix on the estimation of the number of components was investigated as well. For each setting, 100 datasets were generated from the model with a GAR(3) structure. In addition, a different number of mixture components ($k = 2, 3,$ and 4) was considered to evaluate the choice of k . Approximately the correct number of components ($k=3$) was selected for MVTDs in all cases. However, MVNDs tend to overestimate ($k=4$) the number of true components.

3.4 Real data: Gastrointestinal (GI) cancers

The age-standardized death rates of the three most common GI cancers were extracted from the *Our World In Data* website [33]. The information included the death rates (per 100,000 populations) of colon and rectum, stomach, and liver cancers in 186 countries during 1990-2015 (at 5-year intervals). A mixture of MVTDs and MVNDs was fitted with k ranging from 1 to 10. The best sub-model based on the BIC and ICL is (GAR(2), VVV) with $k=7$ in the mixture of MVNDs and (GAR(4), VVV) with the constrained degrees of freedom and $k=6$ in the mixture of MVTDs. The estimated degree of freedom for the mixture of the MVTDs was $\hat{\nu} = 3.33$. We also fitted a finite mixture of skew matrix-variate distributions introduced by *Gallaugher* and *McNicholas* [15] to GI data. These matrix-variate distributions are skew-t, generalized hyperbolic, variance-gamma, and normal inverse Gaussian distributions that we did not consider eigenvalue and the modified Cholesky decompositions for the *between* and *within* covariance matrices for those, respectively. Because of the huge number of parameters, any of these finite mixture had not been converge.

Further, stomach and liver cancer death rates in some countries were extremely higher compared to other countries. Thus, the mixture of the MVND model provided an additional cluster to allow outliers. For more details, the number of countries in each cluster and the maps of the included countries in each cluster of MVTD model is presented in Figure 1. The obtained cluster labels in the two models were the same for 132 countries.

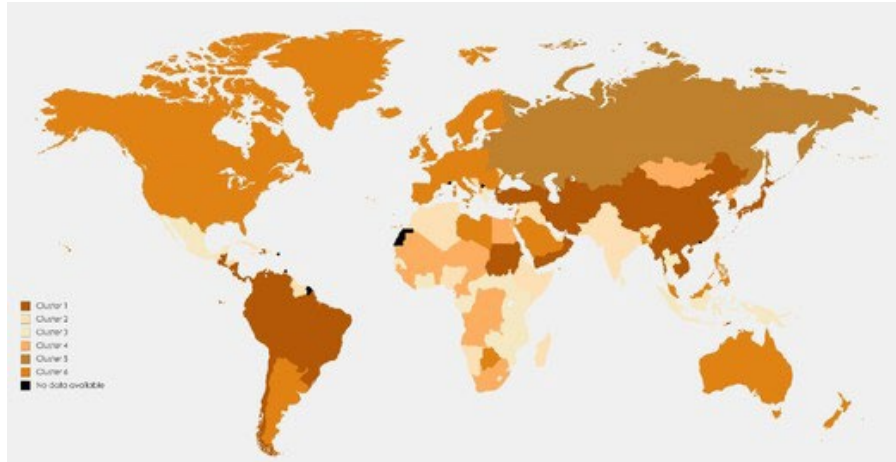


Figure 1: Trend of countries based on the death rates of the three common GI cancer resulted from the mixture of the MVTDS

Conclusion

Based on γ in the mixture models of MVTD and MVND, no differences were observed between the estimation of \mathbf{M} , $\mathbf{\Omega}$, and \mathbf{T} matrices under different orders of temporal covariance structures in each model in simulation studies. Further, these values were similar in both models. On the other hand, the estimation of matrix \mathbf{T} relies on the misspecification of $\mathbf{\Phi}$. Thus, γ_T s should have the least value compared to lower orders if the order of the incorrect temporal structure is equal to or greater than the correct order of the temporal structure. The estimations of matrix \mathbf{T} and the number of mixture components k are overestimated in MVND models if the datasets have a heavy-tail or outlier observations. The mixture of MVTDS commonly selected the model with the right temporal structure and the right number of mixture components. On the other hand, the time it took to fit a mixture of MVTDS was much longer than it required to fit a mixture of MVNDs, which is a trade-off for more precision.

References

- [1] G. J. McLachlan and D. Peel, *Finite mixture models*. New York (NY): John Wiley & Sons, 2000.
- [2] P. D. McNicholas, *Mixture model-based classification*, 1st ed. Chapman & Hall/CRC, 2016.
- [3] G. J. McLachlan and K. E. Basford, *Mixture models : inference and applications to clustering*, 1st ed. CRC Press, 1987.
- [4] D. Peel and G. J. McLachlan, "Robust mixture modelling using the t distribution," *Stat. Comput.*, vol. 10, no. 4, pp. 339–348, 2000.
- [5] T. I. Lin, "Maximum likelihood estimation for multivariate skew normal mixture models," *J. Multivar. Anal.*, vol. 100, no. 2, pp. 257–265, Feb. 2009.
- [6] T. I. Lin, "Robust mixture modeling using multivariate skew t distributions," *Stat. Comput.*, vol. 20, no. 3, pp. 343–356, Jul. 2010.
- [7] A. O'Hagan, T. B. Murphy, I. C. Gormley, P. D. McNicholas, and D. Karlis, "Clustering with the multivariate normal inverse Gaussian distribution," *Comput. Stat. Data Anal.*, vol. 93, pp. 18–30, Jan. 2016.
- [8] R. P. Browne and P. D. McNicholas, "A mixture of generalized hyperbolic distributions," *Can. J. Stat.*, vol. 43, no. 2, pp. 176–198, Jun. 2015.
- [9] U. J. Dang, R. P. Browne, and P. D. McNicholas, "Mixtures of multivariate power exponential distributions," *Biometrics*, vol. 71, no. 4, pp. 1081–1089, Dec. 2015.
- [10] A. K. Gupta and D. K. Nagar, *Matrix Variate Distributions*, 1st ed. New York: Chapman and Hall/CRC, 2000.
- [11] C. Viroli, "Finite mixtures of matrix normal distributions for classifying three-way data," *Stat. Comput.*, vol. 21, no. 4, pp. 511–522, Oct. 2011.
- [12] L. Anderlucci and C. Viroli, "Covariance pattern mixture models for the analysis of multivariate heterogeneous longitudinal data," *Ann. Appl. Stat.*, vol. 9, no. 2, pp. 777–800, 2015.

- [13] F. Z. Doğru, Y. M. Bulut, and O. Arslan, "Finite mixtures of matrix variate t distributions," *Gazi Univ. J. Sci.*, vol. 29, no. 2, pp. 335–341, Jun. 2016.
- [14] M. P. B. Gallagher and P. D. McNicholas, "A matrix variate skew- t distribution," *Stat*, vol. 6, no. 1, pp. 160–170, Jan. 2017.
- [15] M. P. B. Gallagher and P. D. McNicholas, "Finite mixtures of skewed matrix variate distributions," *Pattern Recognit.*, vol. 80, pp. 83–93, Aug. 2018.
- [16] M. P. B. Gallagher and P. D. McNicholas, "Three skewed matrix variate distributions," *Stat. Probab. Lett.*, vol. 145, pp. 103–109, Feb. 2019.
- [17] P. D. McNicholas, T. B. Murphy, A. F. McDaid, and D. Frost, "Serial and parallel implementations of model-based clustering via parsimonious Gaussian mixture models," *Comput. Stat. Data Anal.*, vol. 54, no. 3, pp. 711–723, Mar. 2010.
- [18] J. Banfield and A. Jeffrey, "Model-Based Gaussian and Non-Gaussian Clustering," *Biometrics*, vol. 49, no. 3, pp. 803–21, Sep. 1993.
- [19] G. Celeux and G. Govaert, "Gaussian parsimonious clustering models," *Pattern Recognit.*, vol. 28, no. 5, pp. 781–93, 1995.
- [20] C. Fraley and A. E. Raftery, "Model-Based Clustering, Discriminant Analysis, and Density Estimation.," *J. Am. Stat. Assoc.*, vol. 97, no. 458, pp. 611–631, Jun. 2002.
- [21] P. D. McNicholas and T. B. Murphy, "Model-based clustering of longitudinal data," *Can. J. Stat.*, vol. 38, no. 1, pp. 153–168, 2010.
- [22] J. L. Andrews and P. D. McNicholas, "Model-based clustering, classification, and discriminant analysis via mixtures of multivariate t -distributions," *Stat. Comput.*, vol. 22, no. 5, pp. 1021–1029, 2012.
- [23] P. D. McNicholas and S. Subedi, "Clustering gene expression time course data using mixtures of multivariate t -distributions," *J. Stat. Plan. Inference*, vol. 142, no. 5, pp. 1114–1127, May 2012.
- [24] I. Vrbik and P. D. McNicholas, "Parsimonious skew mixture models for model-based clustering and classification," *Comput. Stat. Data Anal.*, vol. 71, pp. 196–210, Mar. 2014.
- [25] A. K. Gupta, T. Varga, and T. Bodnar, *Elliptically Contoured Models in Statistics and Portfolio Theory*, 2nd ed. New York: Springer, 2013.
- [26] M. P. B. Gallagher and P. D. McNicholas, "A matrix variate skew- t distribution," *Stat*, vol. 6, no. 1, pp. 160–170, 2017.
- [27] S. Sarkar, X. Zhu, V. Melnykov, and S. Ingrassia, "On parsimonious models for modeling matrix data," *Comput. Stat. Data Anal.*, vol. 142, p. 106822, Feb. 2020.
- [28] M. Pourahmadi, "Joint mean-covariance models with applications to longitudinal data: unconstrained parameterisation," *Biometrika*, vol. 86, no. 3, pp. 677–690, Sep. 1999.
- [29] X.-L. MENG and D. B. RUBIN, "Maximum likelihood estimation via the ECM algorithm: A general framework," *Biometrika*, vol. 80, no. 2, pp. 267–278, Jun. 1993.
- [30] G. Schwarz, "Estimating the Dimension of a Model," *Ann. Stat.*, vol. 6, no. 2, pp. 461–464, Mar. 1978.
- [31] C. Biernacki, G. Celeux, and G. Govaert, "Assessing a mixture model for clustering with the integrated completed likelihood," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 7, pp. 719–725, Jul. 2000.
- [32] L. Anderlucci and C. Viroli, "Supplement to "Covariance pattern mixture models for the analysis of multivariate heterogeneous longitudinal data," 2015.
- [33] M. Roser and H. Ritchie, "Cancer," *Our World in Data*, 2019. [Online]. Available: <https://ourworldindata.org/cancer>. [Accessed: 16-Sep-2019].