# Robust mixture modeling using multivariate skew *t* distributions

## Tsung-I Lin

Department of Applied Mathematics and Institute of Statistics

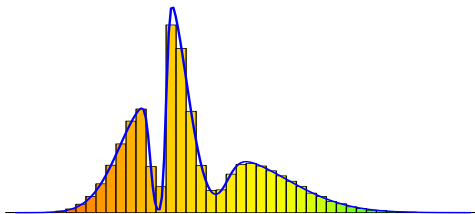National Chung Hsing University, Taiwan

August 24, 2010

National Chung Hsing University

# OUTLINE

# 1. INTRODUCTION

- Finite mixture models have become a useful tool for modeling data that are thought to come from several different groups with varying proportions.



- Lin et al. (2007) proposed a novel (univariate) skew $t$ mixture (STMIX) model, which allows for accommodation of both skewness and thick tails for making robust inferences. Drawback: limited to data with univariate outcomes.

- We propose a multivariate version of the STMIX (MSTMIX) model, composed of a weighed sum of $g$-component multivariate skew $t$ (MST) distributions.

# The multivariate skew $t$ (MST) distribution

- The MST distribution, $\mathbf{Y} \sim \mathcal{S}t_p(\boldsymbol{\xi}, \boldsymbol{\Sigma}, \boldsymbol{\Lambda}, \nu)$, can be represented by

The stochastic representation of skew $t$ distribution

$$\mathbf{Y} = \boldsymbol{\mu} + \frac{\mathbf{Z}}{\sqrt{\tau}}, \ \ \mathbf{Z} \sim \mathcal{S}\mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma}, \boldsymbol{\Lambda}), \ \ \tau \sim \Gamma(\nu/2, \nu/2), \ \ \mathbf{Z} \perp \tau \quad (1)$$

- $\mathbf{Y} \mid \tau \sim \mathcal{S}\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}/\tau, \boldsymbol{\Lambda}/\sqrt{\tau})$

Proposition 1.

If $\tau \sim \Gamma(\alpha, \beta)$, then for any $\boldsymbol{a} \in \mathbb{R}^p$

$$E(\Phi_p(\boldsymbol{a}\sqrt{\tau}|\boldsymbol{\Delta})) = T_p\left(\boldsymbol{a}\sqrt{\frac{\alpha}{\beta}} \ \Big| \ \boldsymbol{\Delta}; 2\alpha\right).$$

- Integrating $\tau$ from the joint density of ($\mathbf{Y}, \tau$) yields

$$\psi(\boldsymbol{y}|\boldsymbol{\xi}, \boldsymbol{\Sigma}, \boldsymbol{\Lambda}, \nu) = 2^p t_p(\boldsymbol{y}|\boldsymbol{\xi}, \boldsymbol{\Omega}, \nu) T_p\left(\boldsymbol{q}\sqrt{\frac{\nu + p}{U + \nu}} \ \Big| \ \boldsymbol{\Delta}; \nu + p\right), \quad (2)$$

where $\boldsymbol{q} = \boldsymbol{\Lambda}\boldsymbol{\Omega}^{-1}(\boldsymbol{y} - \boldsymbol{\xi})$ and $U = (\boldsymbol{y} - \boldsymbol{\xi})^\top \boldsymbol{\Omega}^{-1}(\boldsymbol{y} - \boldsymbol{\xi})$.

$$\boldsymbol{\mu} = \left[ \begin{array}{c} 0 \\ 0 \end{array} \right], \quad \boldsymbol{\Sigma} = \left[ \begin{array}{cc} 1 & \rho \\ \rho & 1 \end{array} \right], \quad \boldsymbol{\lambda} = \left[ \begin{array}{c} \lambda_1 \\ \lambda_2 \end{array} \right], \quad \nu = 4$$
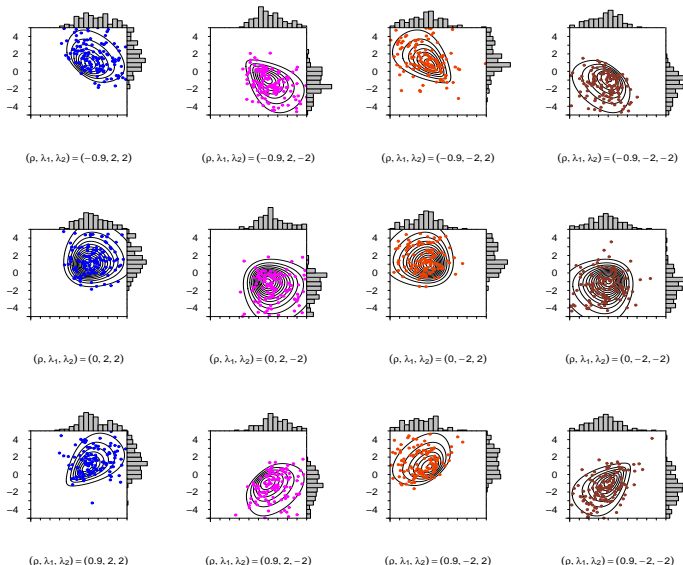


Figure 1: The scatter plots and contours and together with their histograms.

# The MSTMIX model

- The MSTMIX model

$$f(\mathbf{y}_j|\mathbf{\Theta}) = \sum_{i=1}^{g} w_i \psi(\mathbf{y}_j \mid \boldsymbol{\xi}_i, \mathbf{\Sigma}_i, \mathbf{\Lambda}_i, \nu_i), \tag{3}$$

where $\psi(\mathbf{y}_j|\boldsymbol{\xi}_i, \mathbf{\Sigma}_i, \mathbf{\Lambda}_i, \nu_i)$ represents the MST density, and $w_i$'s are the mixing probabilities satisfying $\sum_{i=1}^{g} w_i = 1$.

- Introduce allocation variables $\mathbf{Z}_j = (Z_{1j}, \ldots, Z_{gj})^\top$, $j = 1, \ldots, n$, whose values are a set of binary variables with

$$Z_{ij} = \begin{cases} 1 & \text{if } \mathbf{Y}_j \text{ belongs to group } i, \\ 0 & \text{otherwise}, \end{cases}$$

and satisfying $\sum_{i=1}^{g} Z_{ij} = 1$. Denoted by

$$\mathbf{Z}_j \sim \mathcal{M}(1; w_1, \ldots, w_g).$$

- A hierarchical representation of (3) is

$$
\begin{aligned}
\mathbf{Y}_j \mid (\boldsymbol{\gamma}_j, \tau_j, Z_{ij} = 1) &\sim \mathcal{N}_p(\boldsymbol{\xi}_i + \boldsymbol{\Lambda}_i \boldsymbol{\gamma}_j, \boldsymbol{\Sigma}_i / \tau_j), \\
\boldsymbol{\gamma}_j \mid (\tau_j, Z_{ij} = 1) &\sim \mathcal{HN}_p(\mathbf{0}, \boldsymbol{I}_p / \tau_j), \\
\tau_j \mid (Z_{ij} = 1) &\sim \Gamma(\nu_i/2, \nu_i/2), \\
\mathbf{Z}_j &\sim \mathcal{M}(1; w_1, \ldots, w_g).
\end{aligned}
\tag{4}
$$

- The complete data log-likelihood function of $\boldsymbol{\Theta}$ is

$$
\begin{aligned}
&\ell_c(\boldsymbol{\Theta} \mid \boldsymbol{y}, \boldsymbol{\gamma}, \boldsymbol{\tau}, \mathbf{Z}) \\
&= \sum_{i=1}^{g} \sum_{j=1}^{n} Z_{ij} \Bigg\{ \log(w_i) + \frac{\nu_i}{2} \log\left(\frac{\nu_i}{2}\right) - \log\Gamma\left(\frac{\nu_i}{2}\right) - \frac{1}{2} \log|\boldsymbol{\Sigma}_i| \\
&\quad + \left(\frac{\nu_i}{2} + p - 1\right) \log \tau_j - \frac{\tau_j}{2} \Big( (\boldsymbol{y}_j - \boldsymbol{\xi}_i - \boldsymbol{\Lambda}_i \boldsymbol{\gamma}_j)^\top \boldsymbol{\Sigma}_i^{-1} (\boldsymbol{y}_j - \boldsymbol{\xi}_i - \boldsymbol{\Lambda}_i \boldsymbol{\gamma}_j) \\
&\quad + \nu_i + \boldsymbol{\gamma}_j^\top \boldsymbol{\gamma}_j \Big) \Bigg\}.
\end{aligned}
$$

# Computational aspects of parameter estimation

- The $Q$ function is

$$Q(\boldsymbol{\Theta}|\hat{\boldsymbol{\Theta}}^{(k)}) = E(\ell_c(\boldsymbol{\Theta}|\boldsymbol{y}, \gamma, \tau, \boldsymbol{Z})|\boldsymbol{y}, \hat{\boldsymbol{\Theta}}^{(k)}).$$

- In the MCEM-based algorithm, $Q$-function can be approximated by

$$\hat{Q}(\boldsymbol{\Theta}|\hat{\boldsymbol{\Theta}}^{(k)}) = \frac{1}{M}\sum_{m=1}^{M} \ell_c(\boldsymbol{\Theta} \mid \boldsymbol{y}, \hat{\gamma}_{[m]}^{*(k)}, \hat{\tau}_{[m]}^{*(k)}, \boldsymbol{Z}), \tag{5}$$

where $\hat{\gamma}_{[m]}^{*(k)} = \{\hat{\gamma}_{ij,m}^{*(k)}\}$ and $\hat{\tau}_{[m]}^{*(k)} = \{\hat{\tau}_{ij,m}^{*(k)}\}$ are independently generated by
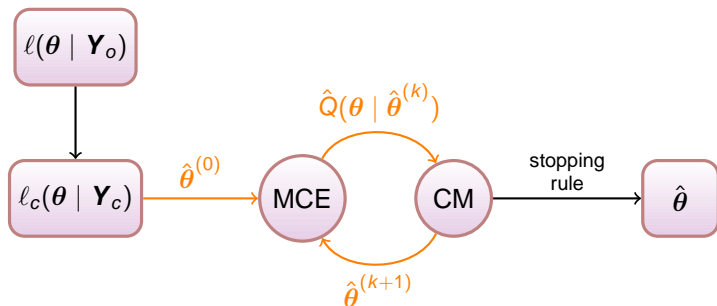
1. $\hat{\gamma}_{ij,m}^{(k+1)}|(\boldsymbol{y}_j, Z_{ij} = 1) \sim \mathcal{T}t_p\Big(\hat{\boldsymbol{q}}_{ij}^{(k)}, \frac{\hat{U}_{ij}^{(k)}+\hat{\nu}_i^{(k)}}{p+\hat{\nu}_i^{(k)}}\hat{\boldsymbol{\Delta}}_i^{(k)}, \hat{\nu}_i^{(k)} + p; \ \mathbb{R}_+^p\Big).$

2. $\hat{\tau}_{ij,m}^{(k+1)}|(\hat{\gamma}_{ij,m}^{(k+1)}, \boldsymbol{y}_j, Z_{ij} = 1)$
   $\sim \Gamma\Big(\frac{\hat{\nu}_i^{(k)} + 2p}{2}, \frac{(\hat{\gamma}_{ij,m}^{(k+1)} - \hat{\boldsymbol{q}}_{ij}^{(k)})^\top \hat{\boldsymbol{\Delta}}_i^{(k)-1}(\hat{\gamma}_{ij,m}^{(k+1)} - \hat{\boldsymbol{q}}_{ij}^{(k)}) + \hat{U}_{ij}^{(k)} + \hat{\nu}_i^{(k)}}{2}\Big).$

# The MCECM algorithm



| $\arg\max Q$ | fix | |
|---|---|---|
| $\boldsymbol{\theta}_1$ | $\hat{\boldsymbol{\theta}}_2^{(k)}$ , | $\hat{\boldsymbol{\theta}}_3^{(k)}$ |
| $\boldsymbol{\theta}_2$ | $\hat{\boldsymbol{\theta}}_1^{(k+1)}$ , | $\hat{\boldsymbol{\theta}}_3^{(k)}$ |
| $\boldsymbol{\theta}_3$ | $\hat{\boldsymbol{\theta}}_1^{(k+1)}$ , | $\hat{\boldsymbol{\theta}}_2^{(k+1)}$ |

**CM-steps:**

$$\hat{w}_i^{(k+1)} = n^{-1} \sum_{j=1}^n \hat{z}_{ij}^{(k)}$$

$$\hat{\boldsymbol{\xi}}_i^{(k+1)} = \frac{\sum_{j=1}^n \hat{\tau}_{ij}^{(k)} \boldsymbol{y}_j - \hat{\boldsymbol{\Lambda}}_i^{(k)} \sum_{j=1}^n \hat{\boldsymbol{\eta}}_{ij}^{(k)}}{\sum_{j=1}^n \hat{\tau}_{ij}^{(k)}}$$

$$\hat{\boldsymbol{\Lambda}}_i^{(k+1)} = \text{diag}\left\{ \left(\hat{\boldsymbol{\Sigma}}_i^{(k)^{-1}} \odot \hat{\boldsymbol{B}}_{1i}^{(k)}\right)^{-1} \left(\hat{\boldsymbol{\Sigma}}_i^{(k)^{-1}} \odot \hat{\boldsymbol{B}}_{2i}^{(k)}\right) \boldsymbol{1}_p \right\}$$

$$\hat{\boldsymbol{\Sigma}}_i^{(k+1)} = \frac{1}{\sum_{j=1}^n \hat{z}_{ij}^{(k)}} \Big( \sum_{j=1}^n \hat{\tau}_{ij}^{(k)} (\boldsymbol{y}_j - \hat{\boldsymbol{\xi}}_i^{(k+1)})(\boldsymbol{y}_j - \hat{\boldsymbol{\xi}}_i^{(k+1)})^\top$$
$$+ \hat{\boldsymbol{\Lambda}}_i^{(k+1)} \hat{\boldsymbol{B}}_{1i}^{(k)} \hat{\boldsymbol{\Lambda}}_i^{(k+1)} - \hat{\boldsymbol{\Lambda}}_i^{(k+1)} \hat{\boldsymbol{B}}_{2i}^{(k)} - \hat{\boldsymbol{B}}_{2i}^{(k)^\top} \hat{\boldsymbol{\Lambda}}_i^{(k+1)} \Big)$$

- Obtain $\hat{\nu}_i^{(k+1)}$ as the solution of
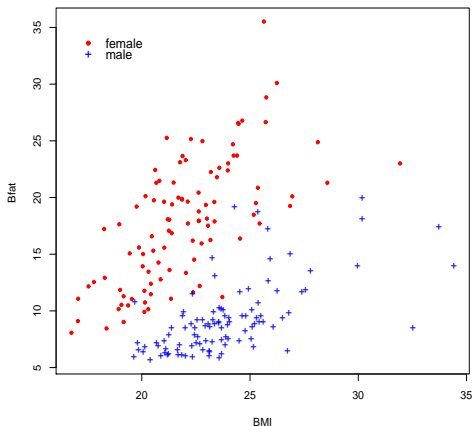
$$\log\left(\frac{\nu_i}{2}\right) + 1 - \text{DG}\left(\frac{\nu_i}{2}\right) + \frac{1}{\sum_{j=1}^n \hat{z}_{ij}^{(k)}} \sum_{j=1}^n (\hat{\kappa}_{ij}^{(k)} - \hat{\tau}_{ij}^{(k)}) = 0.$$

- If the dfs are assumed to be identical, update $\hat{\nu}^{(k)}$ by

$$\hat{\nu}^{(k+1)} = \underset{\nu}{\text{argmax}} \sum_{j=1}^n \log\left( \sum_{i=1}^g \hat{w}_i^{(k+1)} \psi(\boldsymbol{y}_j \mid \hat{\boldsymbol{\xi}}_i^{(k+1)}, \hat{\boldsymbol{\Sigma}}_i^{(k+1)}, \hat{\boldsymbol{\Lambda}}_i^{(k+1)}, \nu) \right).$$

# The Australian Institute of Sport (AIS) data

- Data : The AIS data taken by Cook and Weisberg (1994).
- There are 202 athletes which include 100 females and 102 males.
- Variables : BMI (Body mass index; $kg/m^2$) and Bfat (Body fat percentage).

A two-component MSTMIX model can be written as

$$f(\mathbf{y}_j|\boldsymbol{\Theta}) = wf(\mathbf{y}_j|\boldsymbol{\xi}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\Lambda}_1, \nu_1) + (1-w)f(\mathbf{y}_j|\boldsymbol{\xi}_2, \boldsymbol{\Sigma}_2, \boldsymbol{\Lambda}_2, \nu_2),$$

where

$$\boldsymbol{\xi}_i = (\xi_{i1}, \xi_{i2})^\top, \quad \boldsymbol{\Sigma}_i = \left[ \begin{array}{cc} \sigma_{i,11} & \sigma_{i,12} \\ \sigma_{i,12} & \sigma_{i,22} \end{array} \right] \text{ and } \boldsymbol{\Lambda}_i = \left[ \begin{array}{cc} \lambda_{i,11} & 0 \\ 0 & \lambda_{i,22} \end{array} \right].$$
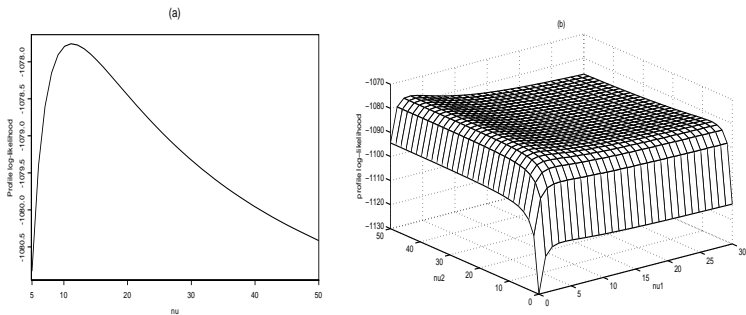


Figure 2: Plot of the profile log-likelihood for $\nu_1$ and $\nu_2$ with a two component MSTMIX model with (a) $\nu_1 = \nu_2 = \nu$ (b) $\nu_1 \neq \nu_2$. ($\hat{\nu}_1 = 4.2$, $\hat{\nu}_2 = 44.1$)

Table 1:Summary results from fitting various mixture models on the AIS data.

| $\Theta$ | MVNMIX | | MVTMIX | | MSNMIX | | MSTMIX | |
|---|---|---|---|---|---|---|---|---|
| | mle | se | mle | se | mle | se | mle | se |
| $w$ | 0.349 | 0.044 | 0.447 | 0.058 | 0.451 | 0.064 | 0.474 | 0.065 |
| $\xi_{11}$ | 23.109 | 0.232 | 23.373 | 2.084 | 21.998 | 2.420 | 21.676 | 0.277 |
| $\xi_{12}$ | 7.959 | 0.203 | 8.320 | 1.428 | 5.898 | 0.141 | 5.947 | 0.057 |
| $\xi_{21}$ | 22.874 | 0.393 | 22.049 | 0.269 | 19.319 | 0.382 | 19.279 | 0.345 |
| $\xi_{22}$ | 16.477 | 0.697 | 17.321 | 0.579 | 13.926 | 1.726 | 17.134 | 1.139 |
| $\sigma_{1,11}$ | 2.878 | 0.700 | 3.791 | 0.873 | 3.178 | 2.988 | 2.730 | 0.392 |
| $\sigma_{1,12}$ | 1.551 | 0.549 | 2.280 | 0.614 | 0.512 | 0.312 | 0.579 | 0.421 |
| $\sigma_{1,22}$ | 2.111 | 0.662 | 3.158 | 0.573 | 0.114 | 0.115 | 0.140 | 0.975 |
| $\sigma_{2,11}$ | 10.971 | 1.468 | 5.606 | 1.098 | 2.765 | 1.055 | 2.420 | 0.533 |
| $\sigma_{2,12}$ | 4.946 | 2.081 | 6.589 | 1.839 | 7.141 | 2.145 | 7.047 | 1.122 |
| $\sigma_{2,22}$ | 32.103 | 4.972 | 24.306 | 5.225 | 20.406 | 9.015 | 23.844 | 0.777 |
| $\lambda_{1,11}$ | — | — | — | — | 1.163 | 3.223 | 1.615 | 0.326 |
| $\lambda_{1,22}$ | — | — | — | — | 3.413 | 0.565 | 3.017 | 0.139 |
| $\lambda_{2,11}$ | — | — | — | — | 4.805 | 0.448 | 4.192 | 1.789 |
| $\lambda_{2,22}$ | — | — | — | — | 4.624 | 1.910 | 0.895 | 6.488 |
| $\nu$ | — | — | 5.820 | 1.646 | — | — | 11.041 | 5.207 |
| $m$ | 11 | | 12 | | 15 | | 16 | |
| $\ell(\hat{\Theta})$ | $-1097.790$ | | $-1093.585$ | | $-1080.647$ | | $-1077.760$ | |
| AIC | 2217.581 | | 2211.170 | | 2191.293 | | 2187.521 | |
| BIC | 2253.972 | | 2250.870 | | 2240.917 | | 2240.453 | |

$\mathrm{AIC} = -2\,\ell(\hat{\Theta}) + 2\,m; \mathrm{BIC} = -2\,\ell(\hat{\Theta}) + m\log(n),\ \ell(\hat{\Theta})$ is the maximized log-likelihood, $m$ is the number of parameters and $n$ is the sample size.
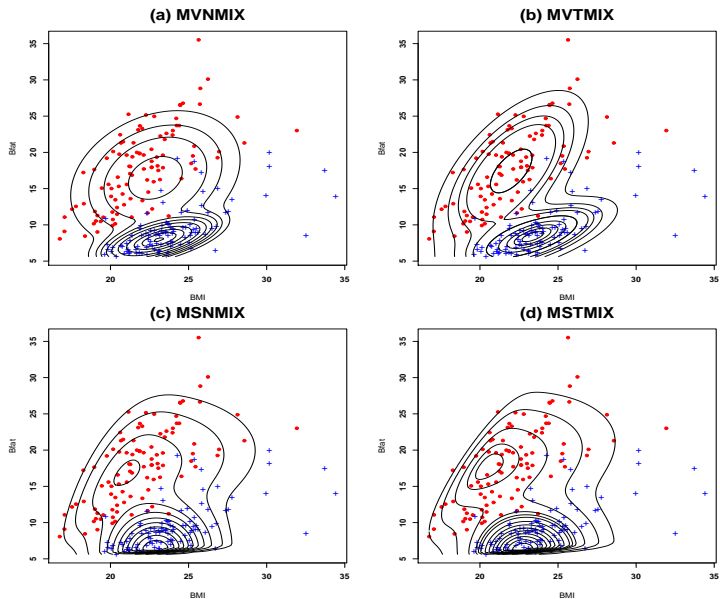
Figure 3: Scatter plot of BMI and Bfat with superimposed contours of two-component various models. The sex are indicated by the female (●) and male (+).

# *Concluding remarks*

- Contributions:

  1. Propose a new robust the MSTMIX model, which offers a great deal of flexibility that accommodates asymmetry and heavy tails simultaneously.

  2. Allow practitioners to analyze heterogeneous multivariate data in a broad variety of considerations.

  3. MCEM-based algorithms are developed for computing ML estimates.

  4. Numerical results show that the MSTMIX model performs reasonably well for the experimental data.