# Model Based Clustering of High-Dimensional Binary Data

Yang Tang        Ryan P. Browne*        Paul D. McNicholas

## Abstract

We propose a mixture of latent trait models with common slope parameters (MCLT) for model-based clustering of high-dimensional binary data, a data type for which few established methods exist. Recent work on clustering of binary data, based on a $d$-dimensional Gaussian latent variable, is extended by incorporating common factor analyzers. Accordingly, our approach facilitates a low-dimensional visual representation of the clusters. We extend the model further by the incorporation of random block effects. The dependencies in each block are taken into account through block-specific parameters that are considered to be random variables. A variational approximation to the likelihood is exploited to derive a fast algorithm for determining the model parameters. Our approach is demonstrated on real and simulated data.

## 1   Introduction

Binary manifest variables are extremely common in behavioural and social sciences research, e.g., individuals may be classified according to whether they take holidays abroad or whether they are satisfied with their lives. In such circumstances, they can be recorded as agreeing or disagreeing with some proposition, or as being capable of doing something or not. Such binary variables are often thought to be indicators of one or more underlying latent variables like, for instance, ability or attitude. Bartholomew and Knott (1999) classify latent variable models into four different classes according to the respective natures of the manifest and latent variables (cf. Table 1). Note that latent trait analysis is termed item response theory (IRT) in the field of educational testing and psychological measurement.

Model-based clustering is a principled statistical approach for clustering, where the data are clustered using some assumed mixture modelling structure. A finite mixture model is a convex combination of a finite number of simple component distributions. Historically, the Gaussian mixture model has dominated the model-based clustering literature (e.g., Wolfe,

*Department of Mathematics &Statistics, University of Guelph, Guelph, Ontario, N1G 2W1, Canada. E-mail: rbrowne@uoguelph.ca.

Table 1: The classification of latent variable methods used by Bartholomew and Knott (1999).

|  |  | Manifest Variables | |
| --- | --- | --- | --- |
|  |  | Metrical | Categorical |
| Latent Variables | Metrical | Factor analysis | Latent trait analysis |
|  | Categorical | Latent profile analysis | Latent class analysis |

1963; Banfield and Raftery, 1993; Celeux and Govaert, 1995; Fraley and Raftery, 2002; McNicholas and Murphy, 2008; Baek et al., 2010). However, very recent model-based clustering work has focused on mixtures of non-elliptical distributions (e.g., Lin, 2010; Lee and McLachlan, 2013; Vrbik and McNicholas, 2012, 2014; Franczak et al., 2014; Murray et al., 2014a,b). Mixture model approaches to data where some or all of the variables are discrete have also been considered (e.g., Hunt and Jorgensen, 1999; McLachlan and Peel, 2000; McLachlan and Chang, 2004).

Model-based approaches for categorical data have received relatively little attention, and recent work on mixtures of latent trait models is summarized in Table 2. Browne and McNicholas (2012) introduce a mixture of latent variables models for the model-based clustering of data with mixed type, and a data set with all binary variables fits within their modelling framework as a special case. They use the deterministic annealing approach to estimate the likelihood described in Zhou and Lange (2010). This approach focuses on increasing the chance of finding the global maximum; however, Gauss-Hermite quadrature is required to approximate the likelihood.

Table 2: Model-based clustering work on discrete data

| Author | Response Function | Data Type | Likelihood Estimation |
| --- | --- | --- | --- |
| Browne and McNicholas (2012) | Logit [a] | Binary/Continuous | Deterministic annealing[b] |
| Cagnone and Viroli (2012) | Logit | Binary | Gauss-Hermite quadrature |
| Gollini and Murphy (2013) | Logit | Binary | Variational EM |
| Muthen and Asparouhov (2006) | Probit | Binary/ Ordered Categorical | Numerical integration |
| Vermunt (2007) | Logit | Multilevel Binary/ Ordered Categorical | Numerical integration |

[a]Binary data fits in the model as a special case.
[b]They use the deterministic annealing approach described by Zhou and Lange (2010).

Gollini and Murphy (2013) propose a mixture of latent trait analyzers (MLTA) for model-based clustering of binary data, wherein a categorical latent variable identifies groups of observations and a latent trait is used to accommodate within cluster dependency. They consider a lower bound approximation to the log-likelihood. This approach is easy to implement and converges quickly in comparison with other numerical approximations to the likelihood. However, mixture of latent trait models become highly parameterized when applied to high-dimensional binary data, particularly when the data come from several different groups and the continuous latent variable is high-dimensional.

A mixture of item response models (Muthen and Asparouhov, 2006; Vermunt, 2007) has very similar structure to the latent trait mixture models; however, it is highly parameterized, uses a probit structure, and numerical integration is required to compute the likelihood. Thus, it can be difficult to apply to large heterogeneous data sets in practice. A similar approach has also been discussed by Cagnone and Viroli (2012), who use Gauss-Hermite quadrature to approximate the likelihood. In addition, they assume a semi-parametric distributional form for the latent variables by adding extra parameters to the model.

Multilevel mixture item response models (Vermunt, 2007) can be used to cluster repeatedly sampled binary data. These models focus on univariate traits because of the number of parameters in the model and the use of quadrature methods for the numerical integration of continuous latent variables. Accordingly, multilevel mixture item response models are not suitable for analyzing large data sets with underlying high-dimensional latent trait structure.

For these reasons, we propose two different mixtures of latent traits models with common slope parameters for model-based clustering of binary data: a general model that supposes that the dependence among the response variables within each observation is wholly explained by a $d$ dimensional continuous latent variable in each group, and an exclusive model for repeatedly sampled data that supposes the response function in each group is composed of two continuous latent variables by adding a blocking latent variable. The proposed family of mixture of latent trait models with common slope parameters (MCLT) is a categorical analogue of a mixture of common factor analyzers model (Baek et al., 2010). The MCLT model enables us to reduce the number of free parameters considerably in estimating the slope. Moreover, it facilitates a low-dimensional visual representation of the clusters with posterior means of the continuous latent variables corresponding to the observed data. The model with a blocking latent variable can potentially reduce known variability of repeatedly sampled data among groups; accordingly, we can be more accurate about group identification.

In the mixture of latent traits model, the likelihood function involves an integral that is intractable. In this work, we propose using a variational approximation of the likelihood, as proposed by Jaakkola and Jordan (2000), Tipping (1999) and Attias (2000), considered a latent variable density model. For a fixed set of values for the variational parameters, the transformed problem has a closed-form solution, providing a lower bound approximation to the log-likelihood. The variational parameters are optimized in a separate step.

The general model is demonstrated on a U.S. Congressional Voting data set (Bache and Lichman, 2013) and the model for clustered data is applied to a data set describing the sensory properties of orange juice (Lee et al., 2013). We compare our approach to the MLTA approach proposed by Gollini and Murphy (2013).

The remainder of this paper is organized as follows. In Section 2, we propose a mixture of latent trait analyzers model with common slope parameters. The data simulations are presented in Section 3. Our approach is then applied to two real data sets (Section 4), and we conclude with a summary and suggestions for future work (Section 5).

# 2 Mixture of Latent Trait Models with Common Slope Parameters

## 2.1 Overview

The MCLT approach restricts the MLTA model by assuming that all latent traits have a set of common slope parameters $\boldsymbol{W} = (\boldsymbol{w}_1, \ldots, \boldsymbol{w}_m)$, for $M$ binary response variables and a $d$-dimensional continuous latent variable $\boldsymbol{Y}$ comes from $g$ different components, where $\boldsymbol{Y}_{ng} \sim \text{MVN}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$. Thus, the MCLT model is a mixture model for binary data that reduces the number of parameters to a manageable size; still, each latent trait has a different effect in each group. It also facilitates low-dimensional visual representation of components with posterior means of the continuous latent variables corresponding to the observed data.

Similar to MLTA model, we assume that each observation $\boldsymbol{x}_n$ $(n = 1, \ldots, N)$ comes from one of the $G$ components and we use $\boldsymbol{z}_n = (z_{n1}, \ldots, z_{nG})$ to identify the group membership, where $z_{ng} = 1$ if observation $n$ is in component $G$ and $z_{ng} = 0$ otherwise. We assume that the conditional distribution of $\boldsymbol{x}_n$ in group $g$ is a latent trait model. Therefore, the MCLT model takes the form,

$$p(\boldsymbol{x}_n) = \sum_{g=1}^{G} \eta_g p(\boldsymbol{x}_n | z_{ng} = 1) = \sum_{g=1}^{G} \eta_g \int_{\boldsymbol{\mathcal{Y}}_{ng}} p(\boldsymbol{x}_n | \boldsymbol{y}_{ng}, z_{ng} = 1) p(\boldsymbol{y}_{ng}) d\boldsymbol{y}_{ng}, \tag{1}$$

where

$$p(\boldsymbol{x}_n | \boldsymbol{y}_{ng}, z_{ng} = 1) = \prod_{m=1}^{M} [\pi_{mg}(\boldsymbol{y}_{ng})]^{x_{nm}} [1 - \pi_{mg}(\boldsymbol{y}_{ng})]^{1-x_{nm}},$$

and the response function for each categorical variable in each group is

$$\pi_{mg}(\boldsymbol{y}_{ng}) = p(x_{nm} = 1 | \boldsymbol{y}_{ng}, z_{ng} = 1) = \frac{1}{1 + \exp\{-\boldsymbol{w}_m' \boldsymbol{y}_{ng}\}}, \tag{2}$$

where $\boldsymbol{w}_m$ is the common model parameter and the latent variable $\boldsymbol{Y}_{ng} \sim \text{MVN}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$.

The complete-data log-likelihood is then given by

$$l = \sum_{n=1}^{N} \log \left[ \sum_{g=1}^{G} \eta_g \int_{\boldsymbol{\mathcal{Y}}_{ng}} \prod_{m=1}^{M} p(x_{nm} | \boldsymbol{y}_{ng}, z_{ng} = 1) p(\boldsymbol{y}_{ng}) d\boldsymbol{y}_{ng} \right]. \tag{3}$$

Therefore, the model is a finite mixture model in which the $g$th component latent variable $\boldsymbol{Y}_{ng}$ is $\text{MVN}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$ and the mixing proportions are $\eta_1, \eta_2, \ldots, \eta_G$.

## 2.2 MCLT with Block Effect

A specific model with block effect can be used for the analysis of clustered data. Clustered data arise, for example, in research designs where a sample of clusters is repeatedly assessed or in educational settings where pupils are clustered within schools. The outcomes stemming from the same cluster tend to be more homogeneous than outcomes stemming from different clusters; accordingly, the outcomes within a cluster are likely to be correlated. These dependencies are taken into account via a response function with a blocking latent variable.

Suppose that each observation $\boldsymbol{x}_{ij}$ is the $j$th observed outcome of cluster $i$ ($i = 1, \ldots, I$; $j = 1, \ldots, J$), and the cluster-specific parameters $s_{ij}$ are assumed to explain all dependencies that are due to inter-cluster variability. Thus, the response function for each group is given by

$$\pi_{mg}(\boldsymbol{y}_{ijg}, s_{ij}) = p(x_{ijm} = 1 | \boldsymbol{y}_{ijg}, s_{ij}, z_{ijg} = 1) = \frac{1}{1 + \exp\{-(\boldsymbol{w}_m' \boldsymbol{y}_{ijg} + \beta_m s_{ij})\}},$$

where $\boldsymbol{w}_m$ and $\beta_m$ are the model parameters. In addition, it is assumed that the blocking latent variable $S_{ij} \sim \mathrm{N}(b_i, \sigma_i^2)$. The model follows naturally:

$$p(\boldsymbol{x}_{ij}) = \sum_{g=1}^{G} \eta_g p(\boldsymbol{x}_{ij} | z_{ijg} = 1) = \sum_{g=1}^{G} \eta_g \int_{\mathbb{R}} \int_{\boldsymbol{\mathcal{Y}}_{ijg}} p(\boldsymbol{x}_{ij} | \boldsymbol{y}_{ijg}, s_{ij}, z_{ijg} = 1) p(\boldsymbol{y}_{ijg}) p(s_{ij}) d\boldsymbol{y}_{ijg} ds_{ij},$$

and the log likelihood can be written

$$l = \sum_{i=1}^{I} \sum_{j=1}^{J} \log \left[ \sum_{g=1}^{G} \eta_g \int_{\mathbb{R}} \int_{\boldsymbol{\mathcal{Y}}_{ijg}} \prod_{m=1}^{M} p(x_{ijm} | \boldsymbol{y}_{ijg}, s_{ij}, z_{ijg} = 1) p(\boldsymbol{y}_{ijg}) p(s_{ij}) d\boldsymbol{y}_{ijg} ds_{ij} \right].$$

The MCLT model with block effect is closely related to the multilevel mixture item response models (Vermunt, 2007; Ng et al., 2006). Vermunt (2007) assumes the latent variables at each level can be continuous, discrete, or both. The MCLT model with block effect is one of the special cases: the lower-level latent variables are combinations of discrete and continuous, with continuous random effects at the higher level. However, we focus on a multivariate trait parameter and the use of common slope parameter considerably reduces the number of free parameters in the model. To the best of our knowledge, this is the first work taking a close look at this particular case.

## 2.3 Gaussian Parsimonious Mixture Models

Following Banfield and Raftery (1993) and Celeux and Govaert (1995), we consider a parametrization of the covariance matrices $\boldsymbol{\Sigma}_1, \ldots, \boldsymbol{\Sigma}_g$ of the component densities. The parametrization of the component covariance matrices via eigenvalue decomposition is

$$\boldsymbol{\Sigma}_g = \lambda_g \boldsymbol{Q}_g \boldsymbol{A}_g \boldsymbol{Q}_g',$$

5

where $\lambda_g = |\boldsymbol{\Sigma}_g|^{\frac{1}{d}}$, $\boldsymbol{Q}_g$ is the matrix of eigenvectors of $\boldsymbol{\Sigma}_g$, and $\boldsymbol{A}_g$ is a diagonal matrix, such that $|\boldsymbol{A}_g| = 1$, with the normalized eigenvalues of $\boldsymbol{\Sigma}_g$ on the diagonal in a decreasing order. The parameter $\lambda_g$ determines the volume of the $g$th cluster, $\boldsymbol{Q}_g$ its orientation, and $\boldsymbol{A}_g$ its shape. We write $\boldsymbol{\Sigma}_g = \lambda_g \boldsymbol{B}_g$, where $\boldsymbol{B}_g$ is a diagonal matrix with $|\boldsymbol{B}_g| = 1$. This particular parametrization gives rise to four models (corresponding to component covariance structures $\lambda \boldsymbol{B}$, $\lambda_g \boldsymbol{B}$, $\lambda \boldsymbol{B}_g$, and $\lambda_g \boldsymbol{B}_g$, respectively). By assuming spherical shapes, namely $\boldsymbol{A}_g = \boldsymbol{I}$, another two parsimonious models are available: $\lambda \boldsymbol{I}$ and $\lambda_g \boldsymbol{I}$. Finally, the 14 parameterizations in Table 3 are considered; note that the corresponding Gaussian mixture models make up the GPCM family of Celeux and Govaert (1995).

Table 3: Fourteen parameterizations of $\boldsymbol{\Sigma_g}$ and the associated number of free parameters.

| | $\boldsymbol{\Sigma}_g$ | Vol/Shape/Orientation[a] | Number of free parameters |
|---|---|---|---|
| 1 | $\lambda \boldsymbol{Q} \boldsymbol{A} \boldsymbol{Q}'$ | EEE | $G - 1 + d(M + G) + d(d+1)/2 - d^2$ |
| 2 | $\lambda_g \boldsymbol{Q} \boldsymbol{A} \boldsymbol{Q}'$ | VEE | $G - 1 + d(M + G) + d(d+1)/2 + G - 1 - d^2$ |
| 3 | $\lambda \boldsymbol{Q} \boldsymbol{A}_g \boldsymbol{Q}'$ | EVE | $G - 1 + d(M + G) + d(d+1)/2 + (G-1)(d-1) - d^2$ |
| 4 | $\lambda_g \boldsymbol{Q} \boldsymbol{A}_g \boldsymbol{Q}'$ | VVE | $G - 1 + d(M + G) + d(d+1)/2 + (G-1)d - d^2$ |
| 5 | $\lambda \boldsymbol{Q}_g \boldsymbol{A} \boldsymbol{Q}'_g$ | EEV | $G - 1 + d(M + G) + G(d(d+1)/2) - (G-1)d - d^2$ |
| 6 | $\lambda_g \boldsymbol{Q}_g \boldsymbol{A} \boldsymbol{Q}'_g$ | VEV | $G - 1 + d(M + G) + G(d(d+1)/2) - (G-1)(d-1) - d^2$ |
| 7 | $\lambda \boldsymbol{Q}_g \boldsymbol{A}_g \boldsymbol{Q}'_g$ | EVV | $G - 1 + d(M + G) + G(d(d+1)/2) - (G-1) - d^2$ |
| 8 | $\lambda_g \boldsymbol{Q}_g \boldsymbol{A}_g \boldsymbol{Q}'_g$ | VVV | $G - 1 + d(M + G) + G(d(d+1)/2) - d^2$ |
| 9 | $\lambda \boldsymbol{B}$ | EEI | $G - 1 + d(M + G) + d - d^2$ |
| 10 | $\lambda_g \boldsymbol{B}$ | VEI | $G - 1 + d(M + G) + G + d - 1 - d^2$ |
| 11 | $\lambda \boldsymbol{B}_g$ | EVI | $G - 1 + d(M + G) + Gd - G + 1 - d^2$ |
| 12 | $\lambda_g \boldsymbol{B}_g$ | VVI | $G - 1 + d(M + G) + Gd - d^2$ |
| 13 | $\lambda \boldsymbol{I}$ | EII | $G - 1 + d(M + G) + 1 - d^2$ |
| 14 | $\lambda_g \boldsymbol{I}$ | VII | $G - 1 + d(M + G) + G - d^2$ |

[a] "E" represents "equal" and "V" represents "variable".

In Table 4, we list the number of parameters to be estimated for the MLTA (Gollini and Murphy, 2013), parsimonious MLTA (PMLTA), and MCLT approaches for $M = 50, 100$, $d = 2$, and $G = 2, 5$. For example, when we cluster $M = 100$ dimensional data into $g = 2$ groups using a $d = 2$ dimensional latent variable, the MLTA model requires 599 parameters to be estimated, while MCLT needs at most 207 parameters. Moreover, as the number of components grows from 2 to 5, the number of parameters grows almost twice as large as before, even for the PMLTA model, but the number of parameters for MCLT remains almost the same.

## 2.4 Interpretation of Model Parameters

The interpretation of the model parameter can be exactly as in MLTA and IRT models. In the finite mixture model, $\eta_g$ is the probability of an observation sampling from the group $g$. The characteristics of component $g$ are determined by a common slope $\boldsymbol{w}_m$, and by the hyperparameters of the latent variable $\boldsymbol{Y}_{ng}$. In the geometric interpretation of the multivariate normal distribution, the equidensity contours of a non-singular multivariate normal

Table 4: The number of free parameters in models for three mixture of latent trait models.

| Model | $G$ | $M$ | $d$ | Number of free parameters |
|---|---|---|---|---|
| MLTA | 2 | 50 | 2 | 299 |
| | 5 | 50 | 2 | 749 |
| | 2 | 100 | 2 | 599 |
| | 5 | 100 | 2 | 1499 |
| PMLTA | 2 | 50 | 2 | 200 |
| | 5 | 50 | 2 | 353 |
| | 2 | 100 | 2 | 400 |
| | 5 | 100 | 2 | 703 |
| MCLT[a] | 2 | 50 | 2 | $102^{b}$–$107^{c}$ |
| | 5 | 50 | 2 | 111–125 |
| | 2 | 100 | 2 | 202–207 |
| | 5 | 100 | 2 | 211–225 |

[a]There are 14 different models for each combination of $G$, $d$ and $M$ (Table 3).
[b]The minimum number of free parameters is calculated by using the coviriance structure $\lambda\boldsymbol{I}$.
[c]The maximum number of free parameters is calculated by using the coviriance structure $\lambda_g\boldsymbol{Q}_g\boldsymbol{A}_g\boldsymbol{Q}'_g$.

distribution are ellipsoids centred at the mean $\boldsymbol{\mu}$. The directions of the principal axes of the ellipsoids are given by the eigenvectors of the covariance matrix $\boldsymbol{\Sigma}$. If $\boldsymbol{Y} \sim \mathrm{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then we have $\boldsymbol{Y} \sim \boldsymbol{\mu} + \boldsymbol{\Sigma}^{\frac{1}{2}}\mathrm{N}(\boldsymbol{0}, \boldsymbol{I})$. Thus, the response function in Equation 2 can be written

$$\pi_{mg}(\boldsymbol{\tau}_n) = p(x_{nm} = 1|\boldsymbol{y}_{ng}, z_{ng} = 1) = \frac{1}{1 + \exp\{-(\boldsymbol{w}'_m\boldsymbol{\mu}_g + \boldsymbol{w}'_m\boldsymbol{\Sigma}_g^{\frac{1}{2}}\boldsymbol{\tau}_n)\}}, \tag{4}$$

where $\boldsymbol{Y}_{ng} \sim \mathrm{N}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$ and $\boldsymbol{\tau}_n \sim \mathrm{N}(\boldsymbol{0}, \boldsymbol{I})$.

Because we have $\boldsymbol{\tau}_n \sim \mathrm{N}(\boldsymbol{0}, \boldsymbol{I})$, the value $\pi_{mg}(0)$ can be used to examine the probability that the median individual in group $g$ has a positive response for the variable $m$,

$$\pi_{mg}(0) = p(x_{nm} = 1|\boldsymbol{\tau}_n = 0, z_{ng} = 1) = \frac{1}{1 + \exp\{-\boldsymbol{w}'_m\boldsymbol{\mu}_g\}}. \tag{5}$$

Moreover, the mean $\boldsymbol{\mu}_g$ and covariance matrix $\boldsymbol{\Sigma}_g$ of component $g$ can be used to provide low-dimensional plots of the cluster.

In the MCLT model with block effect, we can write $\pi^*_{mg}(0)$, which is the $\pi_{mg}(0)$ in (5) adjusted by block effect, as

$$\pi^*_{mg}(0) = \frac{1}{I}\left[\sum_{i=1}^{I} \frac{1}{1 + \exp(-(\boldsymbol{w}'_m\boldsymbol{\mu}_g + \beta'_m s_i))}\right],$$

where $I$ is the number of blocks.

## 2.5  Related Models

There are other statistical models that share a lot of common characteristics with our model. The MCLT model can be treated as a categorical version of a mixture of common factor analyzers (MCFA) model (Baek et al., 2010). The MCFA model employs constraints on the $g$ component means and covariance matrices, i.e.,

$$\boldsymbol{\mu}_g = \boldsymbol{A}\boldsymbol{v}_g, \qquad \text{and} \qquad \boldsymbol{\Sigma}_g = \boldsymbol{A}\boldsymbol{\psi}_g\boldsymbol{A}' + \boldsymbol{D}.$$

A common factor loading matrix $\boldsymbol{A}$ is analogous to a common trait parameter in MCLT. The component covariance matrix is analogous to the covariance matrix of the response function's posterior distribution in each group. The component mean is identical to the mean of the response function's posterior distribution. Of course, the mixing proportions take a same role in both models.

Von Davier and Carstensen (2007) consider a mixture Rasch model (Rasch, 1960) that is equivalent to the parsimonious model of Gollini and Murphy (2013). The model is given by

$$P(x_{nm} = 1 | \boldsymbol{q}_m, \boldsymbol{\beta}_n, \sigma_m, z_{ng} = 1) = \frac{\exp(\boldsymbol{q}_m'\boldsymbol{\beta}_n - \sigma_{mg})}{1 + \exp(\boldsymbol{q}_m'\boldsymbol{\beta}_n - \sigma_{mg})},$$

where $\boldsymbol{q}_m$ are variable-specific parameters, $\boldsymbol{\beta}_n$ is the $d$-dimensional ability parameter, and $\sigma_{mg}$ is the difficulty parameter.

Many other versions of mixture models with response functions have been proposed for analyzing binary data, including the mixed latent trait model (Uebersax, 1999), the latent class factor analysis (LCFA) model (Vermunt et al., 2005), and a range of mixture item response models (Bolt et al., 2001; Muthen and Asparouhov, 2006; Vermunt, 2007, 2008). A key difference between our model and other mixture models is that we focus on a mixture of multivariate latent variables, which allows us to provide low-dimensional plots of the clusters. In addition, we implement a variational approximation for parameter estimation of latent trait models that provides a computationally efficient means of model fitting.

## 2.6  Variational Approximation

Jaakkola and Jordan (2000) introduced a variational approximation for the predictive likelihood in a Bayesian logistic regression model and also briefly considered the "dual" problem, which is closely related to the latent trait model. It obtains a closed form approximation to the posterior distribution of the parameters within a Bayesian framework. Their method is based on a second order Taylor series expansion of the logistic function around a point

$$p(\boldsymbol{x}_{nm} = 1 | \boldsymbol{y}_{ng},\, z_{ng} = 1) = \frac{\exp\{\boldsymbol{w}_m'\boldsymbol{y}_{ng}\}}{1 + \exp\{\boldsymbol{w}_m'\boldsymbol{y}_n\}} = (1 + \exp\{-\boldsymbol{w}_m'\boldsymbol{y}_n\})^{-1}.$$

where $\xi_{nmg} \neq 0$ for all $m = 1, ..., M$. Now, the lower bound of each term in the log-likelihood is given by,

$$L(\boldsymbol{\xi}_{ng}) = \log(\tilde{p}(\boldsymbol{x}_n | \boldsymbol{\xi}_{ng}) = \log \left( \int \prod_{m=1}^{M} \tilde{p}(x_{nm} | \boldsymbol{y}_{ng}, z_{ng} = 1, \xi_{nmg}) p(\boldsymbol{y}_{ng}) \, d\boldsymbol{y}_{ng} \right), \quad (6)$$

where

$$\tilde{p}(x_{nm} | \boldsymbol{y}_{ng}, z_{ng} = 1, \xi_{nmg}) = \sigma(\xi_{nmg}) \exp \left\{ \frac{A_{nmg} - \xi_{nmg}}{2} + \lambda(\xi_{nmg})(A_{nmg}^2 - \xi_{nmg}^2) \right\},$$

$$A_{nmg} = (2x_{nm} - 1)(\boldsymbol{w}_m' \boldsymbol{y}_{ng}), \ \lambda(\xi_{nmg}) = \frac{1}{2\xi_{nmg}} \left[ \frac{1}{2} - \sigma(\xi_{nmg}) \right], \ \sigma(\xi_{nmg}) = (1 + \exp\{-\xi_{nmg}\})^{-1}.$$

This approximation is used to obtain a lower bound for the log-likelihood. A variational EM algorithm (Tipping, 1999) can then be used to obtain parameter estimates that maximize this lower bound.

## 2.7 Model Fitting

When fitting the MCLT model, the integral in the log-likelihood (3) is intractable. Here we illustrate how to use a variational EM algorithm to obtain the approximation of the likelihood:

1. E-Step: estimate $z_{ng}^{(t+1)}$ using

$$z_{ng}^{(t+1)} = \frac{\eta_g^{(t)} \exp\{L(\boldsymbol{\xi}_{ng}^{(t)})\}}{\sum_{g=1}^{G} \eta_g'^{(t)} \exp\{L(\boldsymbol{\xi}'_{ng}^{(t)})\}}.$$

2. M-Step: estimate $\eta_g^{(t+1)}$ using

$$\eta_g^{(t+1)} = \frac{1}{N} \sum_{n=1}^{N} z_{ng}^{(t+1)}.$$

3. Estimate the lower bound of log-likelihood via variational parameter $\xi_{nmg}$:

   (a) E-Step: we approximate the latent posterior statistics for $p(\boldsymbol{y}_{ng} | \boldsymbol{x}_n, z_{ng}^{(t+1)} = 1)$ by its variational lower bound $\underline{p}(\boldsymbol{y}_{ng} | \boldsymbol{x}_n, z_{ng}^{(t+1)} = 1, \boldsymbol{\xi}_{ng}^{(t)})$, which is a $N(\boldsymbol{v}_{ng}^{(t+1)}, \boldsymbol{\varphi}_{ng}^{(t+1)})$ density, where

$$(\boldsymbol{\varphi}_{ng}^{-1})^{(t+1)} = (\boldsymbol{\Sigma}_g^{-1})^{(t)} - 2 \sum_{m=1}^{M} \lambda(\xi_{nmg}^{(t)}) \, \boldsymbol{w}_m^{(t)} \boldsymbol{w}'_m^{(t)},$$

$$\boldsymbol{v}_{ng}^{(t+1)} = \boldsymbol{\varphi}_{ng}^{(t+1)} \left[ (\boldsymbol{\Sigma}_g^{-1})^{(t)} \boldsymbol{\mu}_g^{(t)} + \sum_{m=1}^{M} \left( x_{nm} - \frac{1}{2} \right) \boldsymbol{w}_m^{(t)} \right],$$

where $\sigma(\xi_{nmg}) = (1 + \exp\{-\xi_{nmg}\})^{-1}$, $\lambda(\xi_{nmg}) = (\frac{1}{2} - \sigma(\xi_{nmg}))/2\xi_{nmg}$.

9

(b) M-Step: optimize the variational parameter $\xi_{nmg}^{(t+1)}$. Owing to the EM formulation, each update for $\xi_{nmg}$ corresponds to a monotone improvement to the posterior approximation. The update is

$$(\xi_{nmg}^2)^{(t+1)} = \boldsymbol{w'}_m^{(t)} \mathrm{E}[\boldsymbol{y}_{ng}\boldsymbol{y'}_{ng}]\boldsymbol{w}_m^{(t)},$$

where the expectation is taken with respect to $\underline{p}(\boldsymbol{y}_{ng}|\boldsymbol{x}_n, z_{ng}^{(t+1)} = 1, \boldsymbol{\xi}_{ng}^{(t)})$, the variational posterior distribution based on the previous value of $\xi_{nmg}$. Thus, we have $\mathrm{E}[\boldsymbol{y}_{ng}\boldsymbol{y'}_{ng}] = \boldsymbol{\varphi}_{ng}^{(t+1)} + \boldsymbol{v}_{ng}^{(t+1)}\boldsymbol{v'}_{ng}^{(t+1)}$.

(c) Update parameters $\boldsymbol{w}_m$, $\boldsymbol{\mu}_g$, and $\boldsymbol{\Sigma}_g$ based on the posterior distributions corresponding to the observations in the data set:

$$\boldsymbol{\Sigma}_g^{(t+1)} = \frac{1}{n_g}\sum_{n=1}^{N} z_{ng}^{(t+1)}\boldsymbol{\varphi}_{ng}^{(t+1)},$$

$$\boldsymbol{\mu}_g^{(t+1)} = \frac{1}{n_g}\sum_{n=1}^{N} z_{ng}^{(t+1)}\boldsymbol{v}_{ng}^{(t+1)},$$

where $n_g = z_{1g} + \cdots + z_{Ng}$ and

$$\boldsymbol{w}_m^{(t+1)} = -\left[2\sum_{g=1}^{G}\sum_{n=1}^{N} z_{ng}^{(t+1)}\lambda(\xi_{nmg}^{(t+1)})\,(\boldsymbol{\varphi}_{ng}^{(t+1)} + \boldsymbol{v}_{ng}^{(t+1)}\boldsymbol{v'}_{ng}^{(t+1)})\right]^{-1}$$

$$\times \left[\sum_{g=1}^{G}\sum_{n=1}^{N} z_{ng}^{(i+1)}(x_{nm} - \frac{1}{2})\boldsymbol{v}_{ng}^{(i+1)}\right].$$

(d) Obtain the lower bound of the log-likelihood at the expansion point $\xi_{ng}$:

$$L(\boldsymbol{\xi}_{ng}^{(t+1)}) = \sum_{m=1}^{M}\left[\log\sigma(\xi_{nmg}^{(t+1)}) - \frac{\xi_{nmg}^{(t+1)}}{2} - \lambda(\xi_{nmg}^{(t+1)})(\xi_{nmg}^2)^{(t+1)}\right]$$

$$- \frac{\boldsymbol{\mu'}_g^{(t+1)}(\boldsymbol{\Sigma}_g^{-1})^{(t+1)}\boldsymbol{\mu}_g^{(t+1)}}{2} + \frac{1}{2}\log\frac{|\boldsymbol{\varphi}_{ng}^{(t+1)}|}{|\boldsymbol{\Sigma}_g^{(t+1)}|} + \frac{\boldsymbol{v'}_{ng}^{(t+1)}(\boldsymbol{\varphi}_{ng}^{-1})^{(t+1)}\boldsymbol{v}_{ng}^{(t+1)}}{2},$$

and the log-likelihood:

$$l^{(t)} \approx \sum_{n=1}^{N}\log\left[\sum_{g=1}^{G}\eta_g^{(t+1)}\exp\{L(\boldsymbol{\xi}_{ng}^{(t+1)})\}\right].$$

4. Return to Step 1.

(a) The stopping criterion adopted here is a measure of lack of progress when all parameter estimates become stable and no further improvements can be made to the likelihood value.

$$|\theta^{(t+1)} - \theta^{(t)}| < 0.01.$$

(b) In our application to the voting data (Section 4.1), to facilitate comparison with the MLTA models, convergence of our variational EM algorithm is determined as in Gollini and Murphy (2013). They using a criterion based on the Aitken acceleration (Aitken, 1926), stopping the algorithm when

$$|l_A^{(t+1)} - l_A^{(t)}| \leq 0.01,$$

where

$$l_A^{(t+1)} = l^{(t)} + \frac{1}{1 - a^{(t)}}(l^{(t+1)} - l^{(t)})$$

and

$$a^{(t)} = \frac{l^{(t+1)} - l^{(t)}}{l^{(t)} - l^{(t-1)}}.$$

See Böhning et al. (1994) for details.

5. Approximate the log-likelihood by using the Gauss-Hermite quadrature (Bock and Aitkin, 1981).

### 2.7.1 Model Fitting with Block Effect

To fit model with block effect outlined in Section 2.2, we will need to re-derive the necessary expressions:

1. E-Step: estimate $z_{ijg}^{(t+1)}$ with

$$z_{ijg}^{(t+1)} = \frac{\eta_g^{(t)} \exp\{L(\boldsymbol{\xi}_{ijg}^{(t)})\}}{\sum_{g=1}^{G} \eta_g'^{(t)} \exp\{L(\boldsymbol{\xi}'_{ijg}^{(t)})\}}.$$

2. M-Step: estimate $\eta_g^{(t+1)}$ using

$$\eta_g^{(t+1)} = \frac{\sum_{i=1}^{I} \sum_{j=1}^{J} z_{ijg}^{(t+1)}}{N}.$$

3. Estimate the likelihood

11

(a) E-Step: estimate the latent posterior statistics for $\underline{p}(\boldsymbol{y}_{ijg}, s_{ijg}|\boldsymbol{x}_{ij}, z_{ijg}^{(t+1)} = 1, \boldsymbol{\xi}_{ijg}^{(t)})$ which is a $N(\hat{\boldsymbol{v}}_{ijg}^{(t+1)}, \hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)})$ density:

$$(\hat{\boldsymbol{\varphi}}_{ijg}^{-1})^{(t+1)} = (\hat{\boldsymbol{\Sigma}}_{ijg}^{-1})^{(t)} - 2\sum_{m=1}^{M} \lambda(\xi_{ijmg}^{(t)})\, \hat{\boldsymbol{w}}_m^{(t)} \hat{\boldsymbol{w}}'^{(t)}_m,$$

$$\hat{\boldsymbol{v}}_{ijg}^{(t+1)} = \hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)} \left[ (\hat{\boldsymbol{\Sigma}}_{ijg}^{-1})^{(t)} \hat{\boldsymbol{\mu}}_{ijg}^{(t)} + \sum_{m=1}^{M} (x_{ijm} - \frac{1}{2}) \hat{\boldsymbol{w}}_m^{(t)} \right],$$

where $p(\boldsymbol{y}_{ijg}, s_{ijg})$ is a joint distribution of $p(\boldsymbol{y}_{ijg})$ and $p(s_{ijg})$, which is $N(\hat{\boldsymbol{\mu}}_{ijg}^{(t)}, \hat{\boldsymbol{\Sigma}}_{ijg}^{(t)})$ with $\hat{\boldsymbol{\mu}}_{ijg}^{(t)} = (\boldsymbol{\mu}_g^{(t)}, b_i^{(t)})'$ and

$$\hat{\boldsymbol{\Sigma}}_{ijg}^{(t)} = \begin{pmatrix} \boldsymbol{\Sigma}_g^{(t)} & 0 \\ 0 & (\sigma_i^2)^{(t)} \end{pmatrix}.$$

(b) M-Step: optimize the variational parameter $\xi_{ijmg}^{(t+1)}$:

$$(\xi_{ijmg}^2)^{(t+1)} = \hat{\boldsymbol{w}}'^{(t)}_m \left[ \hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)} + \hat{\boldsymbol{v}}_{ijg}^{(t+1)} \hat{\boldsymbol{v}}'^{(t+1)}_{ijg} \right] \hat{\boldsymbol{w}}_m^{(t)}.$$

(c) Update the parameters $\boldsymbol{w}_m$, $\beta_m$, $\boldsymbol{\mu}_g$, $\boldsymbol{\Sigma}_g$, $b_i$, and $\sigma_i^2$:

$$\boldsymbol{\Sigma}_g^{(t+1)} = \frac{1}{n_g} \sum_{i=1}^{I} \sum_{j=1}^{J} z_{ijg}^{(t+1)} \hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)}, \qquad \boldsymbol{\mu}_g^{(t+1)} = \frac{1}{n_g} \sum_{i=1}^{I} \sum_{j=1}^{J} z_{ijg}^{(t+1)} \hat{\boldsymbol{v}}_{ijg}^{(t+1)},$$

$$(\sigma_i^2)^{(t+1)} = \frac{1}{n_i} \sum_{g=1}^{G} \sum_{j=1}^{J} z_{ijg}^{(t+1)} \hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)}, \qquad b_i^{(t+1)} = \frac{1}{n_i} \sum_{g=1}^{G} \sum_{j=1}^{J} z_{ijg}^{(t+1)} \hat{\boldsymbol{v}}_{ijg}^{(t+1)},$$

where $n_g = \sum_{i=1}^{I} \sum_{j=1}^{J} z_{ijg}$, $n_i = \sum_{g=1}^{G} \sum_{j=1}^{J} z_{ijg}$, and

$$\hat{\boldsymbol{w}}_m^{(t+1)} = - \left( 2 \sum_{g=1}^{G} \sum_{i=1}^{I} \sum_{j=1}^{J} z_{ijg}^{(t+1)} \lambda(\xi_{ijmg}^{(t+1)}) \, (\hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)} + \hat{\boldsymbol{v}}_{ijg}^{(t+1)} \hat{\boldsymbol{v}}'^{(t+1)}_{ijg}) \right)^{-1}$$

$$\times \left( \sum_{g=1}^{G} \sum_{i=1}^{I} \sum_{j=1}^{J} z_{ijg}^{(t+1)} (x_{ijm} - \frac{1}{2}) \hat{\boldsymbol{v}}_{ijg}^{(t+1)} \right),$$

where $\hat{\boldsymbol{w}}_m^{(t+1)} = (w_{m1}^{(t+1)}, \ldots, w_{md}^{(t+1)}, \beta_m^{(t+1)})'$.

(d) Obtain the lower bound of the log-likelihood at the expansion point $\xi_{ijg}$:

$$L(\boldsymbol{\xi}_{ijg}^{(t+1)}) = \sum_{m=1}^{M} \left[ \log \sigma(\xi_{ijmg}^{(t+1)}) - \frac{\xi_{ijmg}^{(t+1)}}{2} - \lambda(\xi_{ijmg}^{(t+1)})(\xi_{ijmg}^2)^{(t+1)} \right]$$

$$- \frac{\hat{\boldsymbol{\mu}}'^{(t+1)}_{ijg} (\hat{\boldsymbol{\Sigma}}_{ijg}^{-1})^{(t+1)} \hat{\boldsymbol{\mu}}_{ijg}^{(t+1)}}{2} + \frac{1}{2} \log \frac{|\hat{\boldsymbol{\varphi}}_{ijg}^{(t+1)}|}{|\hat{\boldsymbol{\Sigma}}_{ijg}^{(t+1)}|} + \frac{\hat{\boldsymbol{v}}'^{(t+1)}_{ijg} (\hat{\boldsymbol{\varphi}}_{ijg}^{-1})^{(t+1)} \hat{\boldsymbol{v}}_{ijg}^{(t+1)}}{2},$$

and the log-likelihood:

$$l^{(t)} \approx \sum_{i=1}^{I} \sum_{j=1}^{J} \log \left[ \sum_{g=1}^{G} \eta_g^{(t+1)} \exp\{L(\boldsymbol{\xi}_{ijg}^{(t+1)})\} \right].$$

## 2.8 Model Selection

We use the Bayesian information criterion (BIC; Schwarz, 1978) as a criterion for model selection,

$$\text{BIC} = -2l + k \log N, \tag{7}$$

where $l$ is the maximized log likelihood, $k$ is the number of free parameters to be estimated in the model, and $N$ is the number of the observations. Within the framework of MCLT models, the values of number of components $G$, the dimension of the latent variable $Y$ (i.e., $d$), and the structure of the connivance matrices $\boldsymbol{\Sigma}_g$ need to be determined. Models with lower values of BIC are preferable. The BIC value could be overestimated using the variational approximation of log-likelihood, which is always less than or equal to the true value. For model selection purposes, we calculate maximized log-likelihood using Gauss-Hermite quadrature after convergence is attained.

For high-dimensional binary data, particularly when the number of observations $n$ is not very large relative to their dimension $m$, it is common to have a large number of patterns with small observed frequency. We cannot use a $\chi^2$ test to check the goodness of the model fit. The analysis of the groups in the selected model can be used to interpret the model. The adjusted Rand index (ARI; Rand, 1971; Hubert and Arabie, 1985) can be used to assess the model performance. The ARI is the corrected-for-chance version of the Rand index. The general form is

$$\frac{\text{index} - \text{expected index}}{\text{maximum index} - \text{expected index}},$$

which is bounded above by 1, and has expected value 0. Intuitively, an ARI value of 1 corresponds to perfect agreement, and a value of 0 would be expected under random classification.

## 2.9 Model Identifiability

The identifiability of our model depends on the identifiability of the latent trait part as well as the identifiability of the mixture models. The identifiability of mixture models has been discussed in McLachlan and Peel (2000). Bartholomew and Knott (1999) give a detailed explanation of model identifiability in the latent trait models context.

The slope parameters $\boldsymbol{W}_g$ are only identifiable with $d \times d$ constraints. This is important when determining the number of free parameters in the model (Table 3).

In addition, Gollini and Murphy (2013) mention that model identifiability holds if the observed information matrix is full rank. This can be checked using empirical methods as

possible non-identifiability can be identified through high standard errors of the parameter estimates and inconsistency between the maximized likelihood values from different random starts.

These checks for identifiability are carried out in our simulation studies (Section 3) and empirical examples (Section 4).

## 2.10  Computational Aspects

We initialize the categorical latent variables $\boldsymbol{z}_n$ $(n = 1, \ldots, N)$ by randomly assigning each observation to one of the $G$ groups. The variational parameters $\xi_{nmg}$ $(n = 1, \ldots, N, \; m = 1, \ldots, M, \; g = 1, \ldots, G)$ are initialized to equal 20, which leads the initial approximation to the conditional distribution to 0. The model parameters are initialized by generating random numbers from a $N(0, 1)$ distribution. The prior means of the latent variable $\boldsymbol{Y}_{ng}$ $(n = 1, \ldots, N \; g = 1, \ldots, G)$ are initialized by random generated number from a $N(0, 1)$. We use $d$-dimensional identity matrices as the initial prior covariance matrices of $\boldsymbol{Y}_{ng}$ $(n = 1, \ldots, N \; g = 1, \ldots, G)$. In addition, the prior mean $\boldsymbol{b} = (b_1, \ldots, b_i)$ and the prior variance $\boldsymbol{\sigma}^2 = (\sigma_1^2, \ldots, \sigma_i^2)$ of the blocking latent variable are set by generating random number from a $N(0, 1)$. We start with ten random initializations of the algorithm and select the model with the lowest BIC.

The use of the variational EM algorithm leads us to an exactly solvable EM algorithm of a latent variable density model that guarantees monotone improvement in the approximation to the likelihood. We also find that this procedure converges rapidly, i.e., only a few iterations are needed.

# 3  Simulation Studies

To illustrate the accuracy of the proposed MCLT model, we performed a simulation experiment on a 20-dimensional binary data set (i.e., $M = 20$). Thus a comparison of approaches (MLTA vs. MCLT) can be carried out. The observations are generated from a MCLT model of the form given in Equation 1 with a two-component mixture ($G = 2, \eta_1 = 0.5$). The latent variables are two-dimensional multivariate normal distributions. The first component has mean $\boldsymbol{\mu}_1 = (0, 1)'$, while the second component has mean $\boldsymbol{\mu}_2 = (3, 3)'$. The covariance matrices take the form, $\boldsymbol{\Sigma}_g = \lambda \boldsymbol{B}_g$. We choose sample sizes $n \in \{100, 250, 500\}$, and run 100 simulations for each sample.

Tables 5 and 6 present the value of true model parameters as well as their mean squared errors (MSE) for $n = 100, 250, 500$. The MSEs decrease with increasing sample size $n$.

In Table 7, we present a comparison of two different approaches on ARI from the clustering results for $n = 100, 250, 500$. Each couplet in Table 7 shows the average ARI and its standard error of ARIs from 100 simulations. With the MCLT approach, the average ARI is 0.64 with a standard error 0.008 for sample size as small as 100 on a 20-dimensional binary

Table 5: True values and the MSEs of $\boldsymbol{w}_m$, tabulated against $n$.

| Variable | Parameters | True | $n=100$ | $n=250$ | $n=500$ | Variable | Parameters | True | $n=100$ | $n=250$ | $n=500$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| M1 | $w_{11}$ | -1.0 | 0.09 | 0.07 | 0.05 | M11 | $w_{111}$ | 0.9 | 0.16 | 0.04 | 0.03 |
|  | $w_{12}$ | -0.7 | 0.13 | 0.06 | 0.07 |  | $w_{112}$ | 0.6 | 0.05 | 0.04 | 0.04 |
| M2 | $w_{21}$ | -0.3 | 0.20 | 0.13 | 0.08 | M12 | $w_{121}$ | -0.4 | 0.05 | 0.05 | 0.03 |
|  | $w_{22}$ | 1.0 | 0.65 | 0.16 | 0.06 |  | $w_{122}$ | 1.7 | 0.31 | 0.15 | 0.05 |
| M3 | $w_{31}$ | 0.88 | 0.36 | 0.09 | 0.04 | M13 | $w_{131}$ | 0.9 | 0.34 | 0.06 | 0.01 |
|  | $w_{32}$ | 0 | 0.37 | 0.09 | 0.08 |  | $w_{132}$ | 0.8 | 0.04 | 0.03 | 0.00 |
| M4 | $w_{41}$ | -0.7 | 0.01 | 0.00 | 0.00 | M14 | $w_{141}$ | 1.5 | 0.09 | 0.01 | 0.00 |
|  | $w_{42}$ | 0.4 | 0.06 | 0.04 | 0.04 |  | $w_{142}$ | 0 | 0.09 | 0.09 | 0.04 |
| M5 | $w_{51}$ | 0.6 | 0.04 | 0.02 | 0.01 | M15 | $w_{151}$ | 1.6 | 0.2 | 0.1 | 0.01 |
|  | $w_{52}$ | -0.4 | 0.06 | 0.06 | 0.05 |  | $w_{152}$ | 0.5 | 0.06 | 0.03 | 0.03 |
| M6 | $w_{61}$ | -0.4 | 0.02 | 0.01 | 0.01 | M16 | $w_{161}$ | -0.5 | 0.3 | 0.1 | 0.03 |
|  | $w_{62}$ | 0 | 0.22 | 0.05 | 0.01 |  | $w_{162}$ | -0.7 | 0.10 | 0.02 | 0.00 |
| M7 | $w_{71}$ | 2 | 0.04 | 0.02 | 0.00 | M17 | $w_{171}$ | -0.5 | 0.02 | 0.02 | 0.00 |
|  | $w_{72}$ | 0.4 | 0.16 | 0.26 | 0.01 |  | $w_{172}$ | -0.7 | 0.01 | 0.01 | 0.00 |
| M8 | $w_{81}$ | -0.5 | 0.00 | 0.00 | 0.00 | M18 | $w_{181}$ | -1.0 | 0.01 | 0.00 | 0.00 |
|  | $w_{82}$ | -0.4 | 0.02 | 0.02 | 0.00 |  | $w_{182}$ | 0.6 | 0.01 | 0.01 | 0.00 |
| M9 | $w_{91}$ | -1 | 0.01 | 0.00 | 0.00 | M19 | $w_{191}$ | 0.0 | 0.00 | 0.00 | 0.00 |
|  | $w_{92}$ | -0.7 | 0.02 | 0.01 | 0.00 |  | $w_{192}$ | 2.8 | 0.01 | 0.03 | 0.00 |
| M10 | $w_{101}$ | 0.7 | 0.00 | 0.00 | 0.00 | M20 | $w_{201}$ | -1.5 | 0.02 | 0.00 | 0.00 |
|  | $w_{102}$ | 0.5 | 0.00 | 0.00 | 0.00 |  | $w_{202}$ | -0.9 | 0.02 | 0.01 | 0.00 |

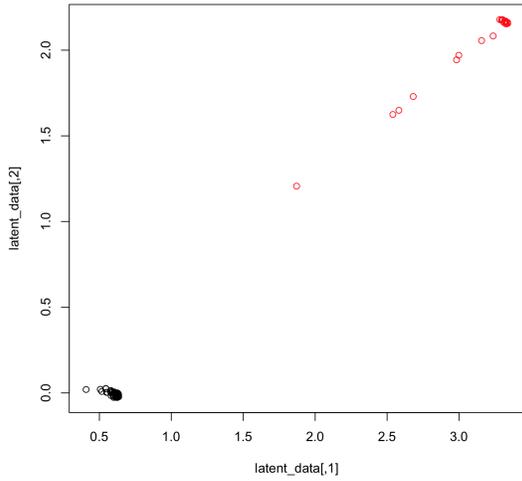Table 6: True values and the MSEs of $\boldsymbol{\mu}_g$, tabulated against $n$.

|  | Parameters | True | $n=100$ | $n=250$ | $n=500$ |
|---|---|---|---|---|---|
| Group 1 | $\mu_{11}$ | 0 | 0.07 | 0.07 | 0.01 |
|  | $\mu_{12}$ | 1 | 0.17 | 0.06 | 0.02 |
| Group 2 | $\mu_{21}$ | 3 | 0.09 | 0.10 | 0.03 |
|  | $\mu_{22}$ | 3 | 0.04 | 0.01 | 0.01 |

data; and a stable clustering result occurs when sample size reaches 250. On the other hand, the average ARI is 0.48 with a standard error of 0.03 for the MLTA approach when $n = 100$; and a stable clustering result only occurs when sample size is 500. The average ARIs using MCLT approach are at least as good as those using MLTA approach for all sample sizes.
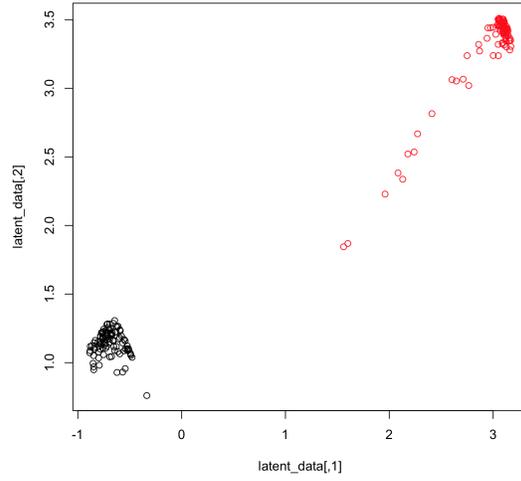
We have also given a plot of the estimated posterior mean for each sample size (Figure 1). These projections are not applicable in MLTA approach as, in its formulation, the latent variables have no cluster-specific discriminatory features.

Table 7: A comparison of two different approaches on ARI and their standard errors, tabulated against $n$.

| Model | $n = 100$ | $n = 250$ | $n = 500$ |
|-------|-----------|-----------|-----------|
| MLTA  | $0.48\,(0.03)$ | $0.54\,(0.02)$ | $0.56\,(0.005)$ |
| MCLT  | $0.65\,(0.008)$ | $0.54\,(0.006)$ | $0.60\,(0.004)$ |



(a) $n = 100$

(b) $n = 250$

(c) $n = 500$

Figure 1: Plots of the estimated posterior mean for different $n$.

# 4 Application

## 4.1 U.S. Congressional Voting

A U.S. congressional voting data set (Bache and Lichman, 2013) has been widely used in the literature (e.g., Gunopulos and Ratanamahatana, 2002; Gollini and Murphy, 2013). This data set includes votes of 435 U.S. House of Representatives congressmen on on sixteen key issues in 1984 with three different type of votes: yes, no, or undecided. The voter's party is labeled as a Democrat or a Republican. The issues voted on are listed in Table 8.

Table 8: The issues that were voted on in the U.S. congressional voting data.

| Item | Issue | Item | Issue |
|------|-------|------|-------|
| 1 | Handicapped Infants | 9 | MX Missile |
| 2 | Water Project Cost-Sharing | 10 | Immigration |
| 3 | Adoption of the Budget Resolution | 11 | Synfuels Corporation Cutback |
| 4 | Physician Fee Freeze | 12 | Education Spending |
| 5 | El Salvador Aid | 13 | Superfund Right to Sue |
| 6 | Religious Groups in Schools | 14 | Crime |
| 7 | Anti-Satellite Test Ban | 15 | Duty- Free Exports |
| 8 | Aid to Nicaraguan 'Contras' | 16 | Export Administration Act/South Africa |

We code each question in two binary variables A and B: the responses for the A variables are coded as 1 = yes/no and 0 = undecided; and B variables are 1 = yes, 0 = no/undecided. The fourteen MCLT models were fitted to these data for $d = 1, 2, \ldots, 5$ and $G = 1, 2, \ldots, 5$. The minimum BIC (Figure 2) occurs at the 2-group, 5-dimensional latent trait model and $\boldsymbol{\Sigma}_g = \lambda \boldsymbol{B}_g$, which is considered as the "best" model. The the BIC value is 9597.

### 4.1.1 A Comparison of approaches: MLTA vs. MCLT

The key statistics on the best models for MLTA, PMLTA, and MCLT are shown in Table 9. It can be seen that the highest ARI value (0.64) is obtained using the MCLT model. Moreover, the MCLT model gives us fewer groups compared to other approaches.

Table 9: Presents a comparison of 3 different approaches.

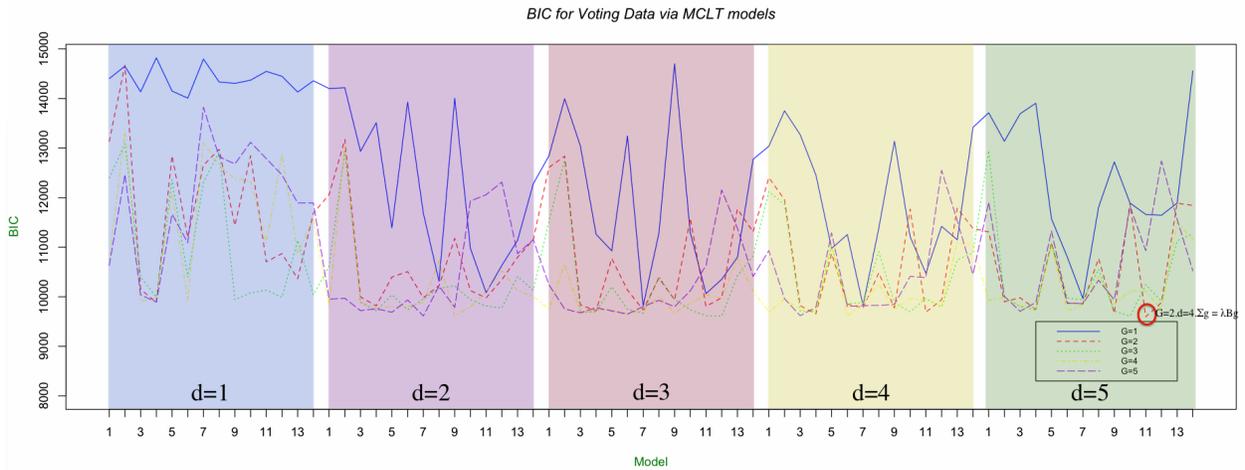| | Model | $G$ | $D$ | BIC | $\boldsymbol{\Sigma}_g$ | ARI |
|---|-------|-----|-----|-----|------------|-----|
| 1 | MLTA | 3 | 1 | 9812 | n/a | 0.42 |
| 2 | MLTA Parsimonious | 4 | 2 | 9681 | n/a | 0.47 |
| 3 | MCLT | 2 | 5 | 9597 | EVI | **0.64** |

Figure 2: BIC values for all 70 different models fitted to the U.S. Congressional voting data for $G = 1, G = 2 \ldots, G = 5$. The order of the models in each dimension ($x$ axis) is as same as in Table 3.

### 4.1.2 Analysis of the Selected MCLT Model

The classification table of the group membership with party membership is presented in Table 10. According to our model selection criteria, BIC=9597 is the minimum BIC with the highest ARI value (0.64) and, therefore, a 2-components and 5-dimensional latent trait model is selected. In comparison with the true party membership, there are only 42 misclassified Congressmen (i.e., 90.3% accuracy) with the "best" model. Group 1 consists mainly of Republican congressman, and Group 2 consists mainly of Democratic congressman. Table 11 shows the median probability $\pi_{mg}(0)$ for each of the groups. The probabilities of a positive response for the A variables (yes/no vs. undecided) for the median individuals in all groups are always high with only one exception in Group 2, for variable number 16, where $\pi_{16\,2}(0) = 0.70$. Thus, the majority of congressmen voted on most issues, but with a slightly lower voting rate in Group 2 on all issues. Due to the high voting rates, most probabilities given for B variables (yes vs. no/undecided) can be interpreted in terms of voting 'yes' versus 'no'.

Table 10: Cross-tabulation of party and predicted classification for our chosen model (EVI, $G = 2$, $d = 5$) for the U.S. Congressional Voting Data.

|  | 1 | 2 |
|---|---|---|
| Republican | 156 | 12 |
| Democrat | 30 | 237 |

It can be observed that the responses for the median individual in Group 1 are opposite to the ones given by the median individual in Group 2 for most issues. The Republican

group (Group 1) tend to give positive responses for the variables 4B, 5B, 6B, 12B, 13B, and 14B. These variables are concerned with the physician fee freeze, El Salvador aid, religious groups in schools, education spending, the superfund right to sue, and crime. The democrat group tend to give positive responses for variables 3B, 7B, 8B, and 9B. These variables are concerned with the adoption of the budget resolutions, the anti-satellite test ban, aid to the Nicaraguan 'Contras', and the MX Missile.

Table 11: Probabilities that the median individual in Group $g$ has a positive response for each of 16 votes in the U.S. Congressional voting data.

| Y/N vs. Undecided | G1 | G2 | Y vs. N/Undecided | G1 | G2 |
|---|---|---|---|---|---|
| 1A | 0.99 | 0.96 | 1B | 0.19 | 0.61 |
| 2A | 0.91 | 0.89 | 2B | 0.50 | 0.41 |
| 3A | 0.98 | 0.97 | 3B | 0.15 | 0.91 |
| 4A | 0.99 | 0.96 | 4B | 0.90 | 0.05 |
| 5A | 0.99 | 0.95 | 5B | 0.98 | 0.10 |
| 6A | 0.99 | 0.96 | 6B | 0.94 | 0.38 |
| 7A | 0.99 | 0.97 | 7B | 0.16 | 0.86 |
| 8A | 0.97 | 0.97 | 8B | 0.08 | 0.93 |
| 9A | 0.99 | 0.93 | 9B | 0.08 | 0.79 |
| 10A | 0.99 | 0.97 | 10B | 0.51 | 0.49 |
| 11A | 0.97 | 0.95 | 11B | 0.21 | 0.43 |
| 12A | 0.95 | 0.92 | 12B | 0.82 | 0.08 |
| 13A | 0.96 | 0.93 | 13B | 0.87 | 0.18 |
| 14A | 0.98 | 0.95 | 14B | 0.97 | 0.27 |
| 15A | 0.95 | 0.93 | 15B | 0.08 | 0.64 |
| 16A | 0.90 | 0.70 | 16B | 0.57 | 0.69 |

We have give a plot of the estimated posterior mean of the best MCLT model with group labels (Figure 3). The two groups are well separated, which can be expected because the error rate of the selected model is quite low (0.093).

## 4.2 Orange Juice Data

A data set describing the sensory properties of orange juice is chosen to illustrate the MCLT model with block effect. The data set contains ten commercially available orange juice (OJ) products. One hundred and twenty consumers were recruited, and the tests were conducted over two weeks in a total of four sessions. The choices within the check-all-that-apply (CATA) questions were presented in alphabetical order during week 1 and in Williams design order (Parker and Williams, 1983) during week 2. In both cases, the attributes were not presented according to sensory modality (appearance, flavour and texture), but in alphabetical order. Therefore, each individual has been accessed 20 times and treated as a block. To the end,

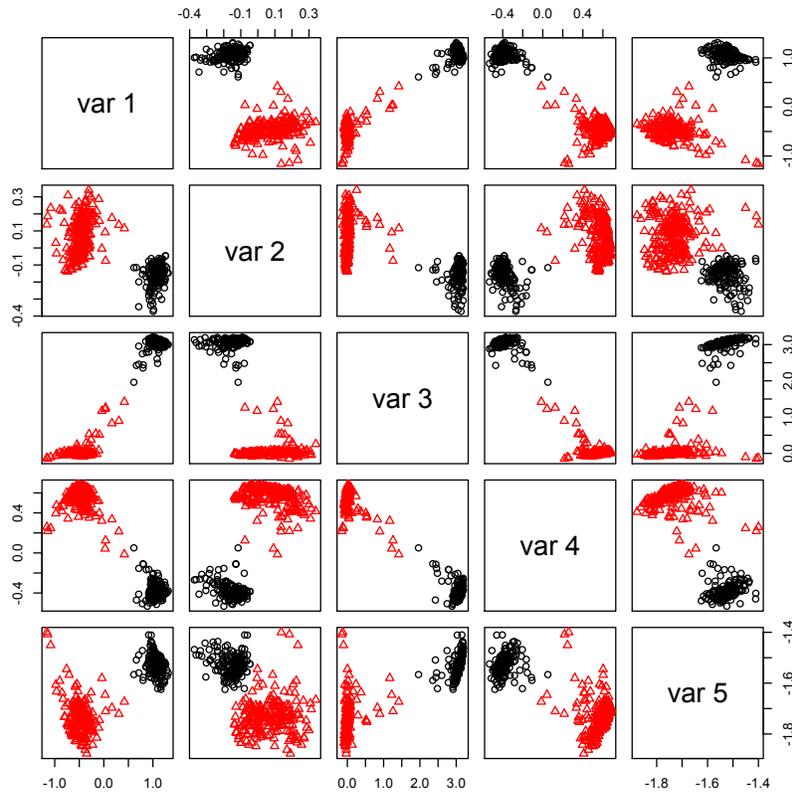Figure 3: Demonstrates the projection of the estimated posterior mean for the **selected MCLT model** via with group labels.

there are 2400 observations, of which 100 are missing. We adopt 40 attributes: 4 in appearance, 27 in flavour, 8 in texture and an indicator for missing observations (Table 12). The study was designed, organized and administered using Compusense® at-hand (Compusense Inc., Guelph, ON, Canada).

Table 12: Attributes for the orange juice data.

| Attribute | Attribute Name | Attribute | Attribute Name |
|---|---|---|---|
| A_1 | A_Cloudy/Turbid | F_18 | F_Other Citrus Flavor |
| A_2 | A_Orange in Color | F_19 | F_Oxidized Flavor |
| A_3 | A_Translucent | F_20 | F_Papery/Cardboard Flavor |
| A_4 | A_Yellow | F_21 | F_Plastic Flavor |
| F_1 | F_Artificial Flavor | F_22 | F_Processed Flavor |
| F_2 | F_Bitter Taste | F_23 | F_Refreshing Flavor |
| F_3 | F_Cheap Taste | F_24 | F_Rotten/Overripe Orange Flavor |
| F_4 | F_Earthy Flavor | F_25 | F_Shelf Stable Flavor |
| F_5 | F_Expensive Flavor | F_26 | F_Strong Flavor |
| F_6 | F_Fresh Orange Flavor | F_27 | F_Weak/Watery Flavor |
| F_7 | F_Fresh Squeezed Flavor | T_1 | T_Astringent/Mouth Drying |
| F_8 | F_From Concentrate Flavor | T_2 | T_Chunky |
| F_9 | F_Green/Unripe Orange Flavor | T_3 | T_Grainy/Chalky |
| F_10 | F_High Acidic/Sour/Tart Taste | T_4 | T_Has a Mouthcoat |
| F_11 | F_High Sweet Taste | T_5 | T_Pulpy |
| F_12 | F_Lemon Flavor | T_6 | T_Smooth |
| F_13 | F_Low Sweet Taste | T_7 | T_Thick |
| F_14 | F_Low Acidic/Sour/Tart Taste | T_8 | T_Thin |
| F_15 | F_Natural Flavor | | |
| F_16 | F_Not From Concentrate Flavor | | |
| F_17 | F_Organic Flavor | | |

We fit MCLT models with block effect to these data for $d = 1, 2, \ldots, 6$ and $G = 1, 2, \ldots, 8$. The minimum BIC (72538) occurs at the 7-group, 4-dimensional latent trait model and $\boldsymbol{\Sigma}_g = \lambda_g \boldsymbol{B}$, which is considered as the "best" model (Table 13).

Table 13: BIC values for model VEI ($\boldsymbol{\Sigma}_g = \lambda_g \boldsymbol{B}$) are listed.

| Dim/Group | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| 1 | 86709 | 85557 | 83375 | 84605 | 86624 | 81715 | 82404 | 78416 |
| 2 | 83358 | 84556 | 82002 | 80160 | 81757 | 81577 | 76281 | 75280 |
| 3 | 81984 | 80535 | 80083 | 80248 | 80614 | 78168 | 75646 | 75014 |
| 4 | 83914 | 80576 | 79178 | 82180 | 77992 | 76571 | **72538** | 78414 |
| 5 | 84556 | 81765 | 80160 | 84350 | 77039 | 77794 | 75009 | 78800 |
| 6 | 85427 | 83044 | 83703 | 85389 | 75376 | 77986 | 75443 | 80535 |

### 4.2.1 Analysis of the Selected MCLT Model

Instead of treating each observation independently, we treat each individual as a block, where each block consists of 20 observations. The classification table of the group membership with product label is presented in Table 14. Group 1 consists mainly of products 4, 6, 7, 10; Group 2 consists mainly products 1, 2, 3, 5, 8; Group 3 consists mainly of the missing observations; Group 4 is a small group consists mainly products 1, 2, 5, 9; Group 5 has only four observations; Group 6 is another small group consists mainly products 4, 7, 10; and Group 7 consists mainly of products 2, 3, 5, 8, 9.

Table 14: Cross-tabulation of predicted classifications versus product label for the best MCLT model (VEI, $G = 7$, $d = 4$) applied to the orange juice data.

| | $G=1$ | $G=2$ | $G=3$ | $G=4$ | $G=5$ | $G=6$ | $G=7$ |
|---|---|---|---|---|---|---|---|
| Missing | 0 | 0 | **100** | 0 | 0 | 0 | 0 |
| P1 | 48 | **103** | 42 | 11 | 1 | 0 | 23 |
| P2 | 25 | **115** | 21 | 14 | 0 | 1 | 55 |
| P3 | 21 | **125** | 21 | 9 | 0 | 4 | 49 |
| P4 | **133** | 48 | 35 | 2 | 0 | 13 | 1 |
| P5 | 16 | **117** | 27 | 19 | 0 | 0 | 54 |
| P6 | **170** | 28 | 22 | 0 | 1 | 6 | 5 |
| P7 | **162** | 35 | 19 | 1 | 0 | 10 | 2 |
| P8 | 12 | **138** | 20 | 7 | 1 | 0 | 53 |
| P9 | 39 | 79 | 27 | 13 | 1 | 1 | 68 |
| P10 | **121** | 67 | 21 | 1 | 0 | 15 | 2 |
| Average Overall Impression | 5.34 | 7.5 | 6.11 | 5.97 | 5.75 | 7.46 | 4.91 |

The groups found in this analysis have similar structures to the ones found using MLTA. By adding the blocking latent variable, we can separate products more accurately. From Table 16 it can be seen that Group 1 consists mainly of products that appear yellow, taste artificial, and have thin texture. The average overall impression score of Group 1 is 5.3/10, which is relatively low among all groups. In contrast, Group 2 consists mainly of products that are orange in colour, fresh in flavour, and pulpy in texture. The average overall impression score of Group 2 is 7.5/10 which is the highest among all groups. Group 7 is characterized by products are thick, taste bitter, and look cloudy to consumers. The average overall impression score of Group 7 is 4.9/10 which is the lowest among all groups. Group 6 is a small group consists of products that are thin but smooth in texture. Group 4 is another small group consists of products have pulpy texture. All missing observations fall into Group 3, and all other observations therein have low probability of positive responses for all attributes (cf. Table 16).

Because there are a large number of response patterns with a very small number of observations (of all $2,241$ observed response patterns, only 50 contain more than one count

and only 7 contain more than two counts), the Pearson's $\chi^2$ test is not applicable. We calculate the overall number of selections (counts) for each attribute across all products to check the goodness of fit. Table 15 shows the observed counts and the expected counts for counts over 500. The table shows that there is a close match between the observed and expected frequencies for most attributes under this model.

Table 15: Observed and expected counts for attributes with 500 or more obsered counts.

|  | Attribute | Observed Counts | Expected Counts |
|---|---|---|---|
| Apperance | A_1 | 777 | 653 |
|  | A_2 | 1364 | 1583 |
|  | A_4 | 924 | 871 |
| Flavour | F_1 | 539 | 557 |
|  | F_2 | 529 | 221 |
|  | F_6 | 757 | 598 |
|  | F_7 | 595 | 597 |
|  | F_8 | 623 | 446 |
|  | F_10 | 563 | 135 |
|  | F_13 | 614 | 114 |
|  | F_15 | 586 | 598 |
|  | F_22 | 512 | 333 |
|  | F_23 | 608 | 598 |
|  | F_26 | 717 | 683 |
| Texture | T_5 | 1128 | 1269 |
|  | T_6 | 895 | 867 |
|  | T_7 | 641 | 926 |
|  | T_8 | 667 | 690 |

We have also given the plot of the estimated posterior mean in selected model via MCLT with group labels (Figure 4). Despite the fact that the posterior mean of Groups 3 and 4 are close together in $d_1$ and $d_2$, all groups are well separated in all dimensions.

# 5    Conclusion

The mixture of latent trait models with common slope parameters gives good clustering performance when applied to high-dimensional binary data. The MCLT model with block effect provides a suitable alternative for clustered data. Our variational EM algorithm gives provided an effective and efficient approach to parameter estimation.

The MCLT model provides a model-based clustering framework for high-dimensional binary data by drawing on ideas from common factor analyzers. The sharing of the slope
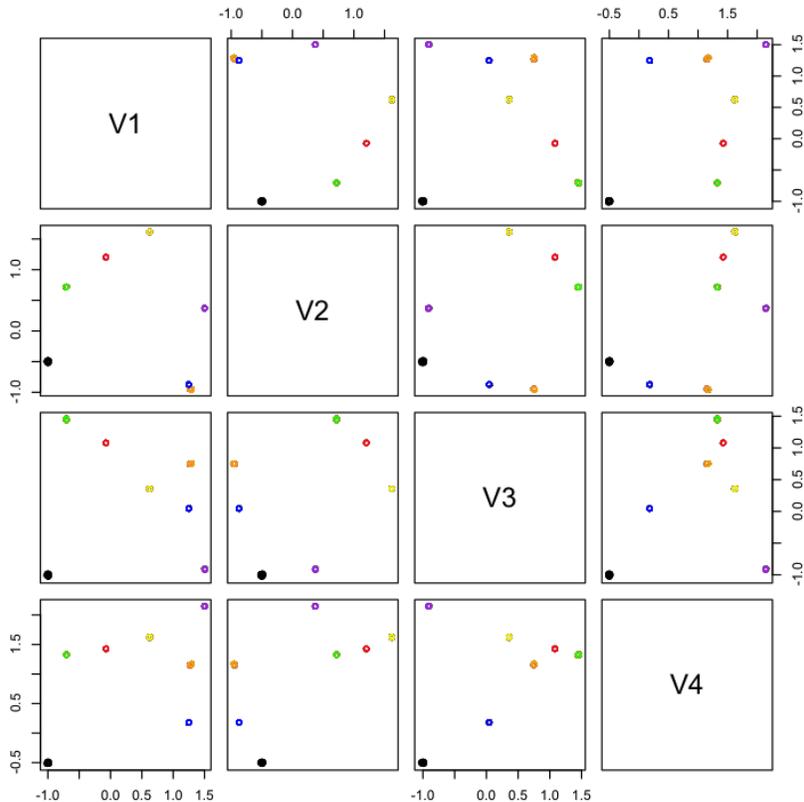
Figure 4: Projects the estimated posterior mean of the best model via MCLT approach with group labels.

parameters enables the model to cluster high-dimensional binary data and to provide low-dimensional plots of the clusters so obtained. The latter plots are given in terms of the (estimated) posterior means of the latent variables. These projections are not applicable in the MLTA approach as, in its formulation, the latent variables have no cluster-specific discriminatory features. The MLTA approach does allow a more general representation of the component covariances and places no restrictions on the component means. However, in this paper, we demonstrate that the MCLT model is useful when the dimension $m$ and the number of clusters $G$ is large. In analogy to the famous GPCM family of mixture models of Celeux and Govaert (1995), c.f. Section 2.3, fourteen covariance structures have been implemented to introduce parsimony. We have presented analyses of two data sets to demonstrate the usefulness of this approach. The model parameters are interpretable and provide a characterization of the within-cluster structure. In our applications herein, we used the BIC to choose the number of clusters $G$, the latent variable dimension $d$, and the covariance decomposition.

In our future work, we wish to investigate other alternatives for repeatedly sampled data. An alternative model for multi-nominal data can be developed using a mixture polytomous logit model.

# Acknowledgements

# References

Aitken, A. C. (1926). On bernoullis numerical solution of algebraic equations. In *Proc. Roy. Soc. Edinburgh*, Volume 46, pp. 289.

Asuncion, A. and D. J. Newman (2007). UCI machine learning repository.

Attias, H. (2000). A Variational Bayesian Framework for Graphical Models. In *In Advances in Neural Information Processing Systems 12*, 209–215

Bache, K. and M. Lichman (2013). UCI machine learning repository.

Baek, J., G. J. McLachlan, and L. K. Flack (2010). Mixtures of factor analyzers with common factor loadings: Applications to the clustering and visualization of high-dimensional data. *IEEE Transactions on Pattern Analysis and Machine Intelligence 32*(7), 1298–1309.

Banfield, J. D. and A. E. Raftery (1993). Model-based gaussian and non-gaussian clustering. *Biometrics*, 803–821.

Bartholomew, D. J. and M. Knott (1999). *Latent Variable Models and Factor Analysis*. Number 7. Edward Arnold.

Bock, D. R. and M. Aitkin (1981). Marginal maximum likelihood estimation of item parameters: Application of an em algorithm. *Psychometrika 46*(4), 443–459.

Böhning, D., E. Dietz, R. Schaub, P. Schlattmann, and B. Lindsay (1994). The distribution of the likelihood ratio for mixtures of densities from the one-parameter exponential family. *Annals of the Institute of Statistical Mathematics 46*, 373–388.

Bolt, D. M., A. S. Cohen, and J. A. Wollack (2001). A mixture item response model for multiple-choice data. *Journal of Educational and Behavioral Statistics 26*(4), 381–409.

Browne, R. P. and P. D. McNicholas (2012). Model-based clustering, classification, and discriminant analysis of data with mixed type. *Journal of Statistical Planning and Inference 142*(11), 2976–2984.

Cagnone, S. and C. Viroli (2012). A factor mixture analysis model for multivariate binary data. *Statistical Modelling 12*(3), 257–277.

Celeux, G. and G. Govaert (1995). Gaussian parsimonious clustering models. *Pattern recognition 28*(5), 781–793.

Fraley, C. and A. E. Raftery (2002). Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association 97*(458), 611–631.

Franczak, B. C., R. P. Browne, and P. D. McNicholas (2014). Mixtures of shifted asymmetric Laplace distributions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. in press.

Gollini, I. and T. B. Murphy (2013). Mixture of latent trait analyzers for model-based clustering of categorical data. *Statistics and Computing*, 1–20.

Gunopulos, D. and C. A. Ratanamahatana (2002). Scaling up the naive bayesian classifier: Using decision trees for feature selection.

Hubert, L. and P. Arabie (1985). Comparing partitions. *Journal of classification 2*(1), 193–218.

Hunt, L. A. and M. A. Jorgensen (1999). Mixture model clustering: a brief introduction to the multimix program. *Australian and New Zealand Journal of Statistics 40*, 153–171.

Jaakkola, T. S. and M. I. Jordan (2000). Bayesian parameter estimation via variational methods. *Statistics and Computing 10*(1), 25–37.

Lee, S. X. and G. J. McLachlan (2013). On mixtures of skew normal and skew t-distributions. *Advances in Data Analysis and Classification 7*(3), 241–266.

Lee, Y., C. Findlay, and J. Meullenet (2013). Experimental consideration for the use of check-all-that-apply questions to describe the sensory properties of orange juices. *International Journal of Food Science & Technology 48*(1), 215–219.

Lin, T.-I. (2010). Robust mixture modeling using multivariate skew t distributions. *Statistics and Computing 20*(3), 343–356.

McLachlan, G. and D. Peel (2000). *Finite Mixture Models*. Wiley.

McLachlan, G. J. and S. U. Chang (2004). Mixture modelling for cluster analysis. *Statistical Methods in Medical Research 13*, 347–361.

McNicholas, P. D. and T. B. Murphy (2008). Parsimonious Gaussian mixture models. *Statistics and Computing 18*(3), 285–296.

Murray, P. M., R. P. Browne, and P. D. McNicholas (2014a). Mixtures of skew-t factor analyzers. *Computational Statistics and Data Analysis*. in press.

Murray, P. M., P. D. McNicholas, and R. P. Browne (2014b). A mixture of common skew-t factor analyzers. *Stat 3*(1), 68–82.

Muthen, B. and Asparouhov (2006). Item response mixture modeling: Application to tobacco dependence criteria. *Addictive Behaviors 31*(6), 1050–1066.

Ng, S., G. McLachlan, K. Wang, L. B. Jones, and S. Ng (2006). A mixture model with random-effects components for clustering correlated gene-expression profiles. *Bioinformatics 22*(14), 1745–1752.

Parker, K. P. and T. W. Williams (1983). Design for testabilitya survey. *Proceedings of the IEEE 71*(1), 98–112.

Rand, W. M. (1971). Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association 66*(336), 846–850.

Rasch, G. (1960). Studies in mathematical psychology: I. probabilistic models for some intelligence and attainment tests.

Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Statist. 6*(2), 461–464.

Tipping, M. E. (1999). Probabilistic visualisation of high-dimensional binary data. *Advances in neural information processing systems*, 592–598.

Uebersax, J. S. (1999). Probit latent class analysis with dichotomous or ordered category measures: conditional independence/dependence models. *Applied Psychological Measurement 23*(4), 283–297.

Vermunt, J. K. (2007). Multilevel mixture item response theory models: an application in education testing. In *Proceedings of the 56th session of the International Statistical Institute*, Lisbon, Portugal, pp. 22–28.

Vermunt, J. K. (2008). Latent class and finite mixture models for multilevel data sets. *Statistical Methods in Medical Research 17*(1), 33–51.

Vermunt, J. K., J. Magidson, and S. I. Inc (2005). Factor analysis with categorical indicators: A comparison between traditional and latent class approaches. *New developments in categorical data analysis for the social and behavioral sciences*, 41–62.

Von Davier, M. and C. H. Carstensen (2007). *Multivariate and mixture distribution Rasch models: Extensions and applications.* Springer Science+ Business Media.

Vrbik, I. and P. D. McNicholas (2012). Analytic calculations for the EM algorithm for multivariate skew-mixture models. *Statistics and Probability Letters 82*(6), 1169–1174.

Vrbik, I. and P. D. McNicholas (2014). Parsimonious skew mixture models for model-based clustering and classification. *Computational Statistics and Data Analysis 71*, 196–210.

Wolfe, J. H. (1963). Object cluster analysis of social areas. Master's thesis, University of California, Berkeley.

Zhou, H. and K. L. Lange (2010). On the bumpy road to the dominant mode. *Scandinavian Journal of Statistics 37*(4), 612–631.

# A    Additional Table

Table 16: Probabilities that the median individual in Group $G$ has a positive response for each attribute in the orange juice data.

| Attribute | Group | $\pi_{mg}(0)$ | Attribute | Group | $\pi_{mg}(0)$ | Attribute | Group | $\pi_{mg}(0)$ |
|---|---|---|---|---|---|---|---|---|
| A_1 | 1 | 0.26 | F_10 | 1 | 0.37 | F_23 | 1 | 0.05 |
|  | 2 | 0.30 |  | 2 | 0.10 |  | 2 | 0.61 |
|  | 3 | 0.27 |  | 3 | 0.16 |  | 3 | 0.06 |
|  | 4 | 0.33 |  | 4 | 0.06 |  | 4 | 0.19 |
|  | 5 | 0.39 |  | 5 | 0.25 |  | 5 | 0.18 |
|  | 6 | 0.10 |  | 6 | 0.10 |  | 6 | 0.54 |
|  | 7 | 0.62 |  | 7 | 0.49 |  | 7 | 0.02 |
| A_2 | 1 | 0.43 | F_11 | 1 | 0.13 | F_24 | 1 | 0.08 |
|  | 2 | 0.72 |  | 2 | 0.22 |  | 2 | 0.03 |
|  | 3 | 0.40 |  | 3 | 0.09 |  | 3 | 0.03 |
|  | 4 | 0.08 |  | 4 | 0.21 |  | 4 | 0.11 |
|  | 5 | 0.67 |  | 5 | 0.18 |  | 5 | 0.10 |
|  | 6 | 0.19 |  | 6 | 0.23 |  | 6 | 0.02 |
|  | 7 | 0.90 |  | 7 | 0.09 |  | 7 | 0.19 |
| A_3 | 1 | 0.14 | F_12 | 1 | 0.12 | F_25 | 1 | 0.07 |
|  | 2 | 0.07 |  | 2 | 0.05 |  | 2 | 0.04 |
|  | 3 | 0.04 |  | 3 | 0.03 |  | 3 | 0.02 |
|  | 4 | 0.06 |  | 4 | 0.05 |  | 4 | 0.04 |
|  | 5 | 0.09 |  | 5 | 0.10 |  | 5 | 0.06 |
|  | 6 | 0.22 |  | 6 | 0.06 |  | 6 | 0.08 |
|  | 7 | 0.03 |  | 7 | 0.10 |  | 7 | 0.03 |
| A_4 | 1 | 0.50 | F_13 | 1 | 0.28 | F_26 | 1 | 0.26 |
|  | 2 | 0.27 |  | 2 | 0.22 |  | 2 | 0.34 |
|  | 3 | 0.48 |  | 3 | 0.23 |  | 3 | 0.18 |
|  | 4 | 0.78 |  | 4 | 0.30 |  | 4 | 0.07 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 5 | 0.27 | | 5 | 0.29 | | 5 | 0.35 |
| | 6 | 0.81 | | 6 | 0.20 | | 6 | 0.15 |
| | 7 | 0.07 | | 7 | 0.32 | | 7 | 0.51 |
| F_1 | 1 | 0.50 | F_14 | 1 | 0.18 | F_27 | 1 | 0.28 |
| | 2 | 0.04 | | 2 | 0.23 | | 2 | 0.03 |
| | 3 | 0.16 | | 3 | 0.17 | | 3 | 0.09 |
| | 4 | 0.15 | | 4 | 0.47 | | 4 | 0.21 |
| | 5 | 0.18 | | 5 | 0.22 | | 5 | 0.12 |
| | 6 | 0.15 | | 6 | 0.26 | | 6 | 0.12 |
| | 7 | 0.30 | | 7 | 0.13 | | 7 | 0.12 |
| F_2 | 1 | 0.39 | F_15 | 1 | 0.03 | T_1 | 1 | 0.27 |
| | 2 | 0.03 | | 2 | 0.60 | | 2 | 0.05 |
| | 3 | 0.12 | | 3 | 0.05 | | 3 | 0.08 |
| | 4 | 0.07 | | 4 | 0.18 | | 4 | 0.05 |
| | 5 | 0.20 | | 5 | 0.20 | | 5 | 0.14 |
| | 6 | 0.04 | | 6 | 0.25 | | 6 | 0.12 |
| | 7 | 0.57 | | 7 | 0.04 | | 7 | 0.18 |
| F_3 | 1 | 0.42 | F_16 | 1 | 0.06 | T_2 | 1 | 0.03 |
| | 2 | 0.02 | | 2 | 0.25 | | 2 | 0.10 |
| | 3 | 0.10 | | 3 | 0.04 | | 3 | 0.03 |
| | 4 | 0.12 | | 4 | 0.09 | | 4 | 0.21 |
| | 5 | 0.13 | | 5 | 0.15 | | 5 | 0.16 |
| | 6 | 0.11 | | 6 | 0.12 | | 6 | 0.01 |
| | 7 | 0.20 | | 7 | 0.07 | | 7 | 0.32 |
| F_4 | 1 | 0.06 | F_17 | 1 | 0.03 | T_3 | 1 | 0.09 |
| | 2 | 0.12 | | 2 | 0.12 | | 2 | 0.04 |
| | 3 | 0.03 | | 3 | 0.01 | | 3 | 0.03 |
| | 4 | 0.07 | | 4 | 0.07 | | 4 | 0.10 |
| | 5 | 0.11 | | 5 | 0.09 | | 5 | 0.10 |
| | 6 | 0.06 | | 6 | 0.05 | | 6 | 0.02 |
| | 7 | 0.08 | | 7 | 0.04 | | 7 | 0.14 |
| F_5 | 1 | 0.02 | F_18 | 1 | 0.23 | T_4 | 1 | 0.22 |
| | 2 | 0.24 | | 2 | 0.08 | | 2 | 0.10 |
| | 3 | 0.01 | | 3 | 0.10 | | 3 | 0.08 |
| | 4 | 0.05 | | 4 | 0.10 | | 4 | 0.05 |
| | 5 | 0.09 | | 5 | 0.17 | | 5 | 0.18 |
| | 6 | 0.09 | | 6 | 0.10 | | 6 | 0.10 |
| | 7 | 0.02 | | 7 | 0.22 | | 7 | 0.24 |
| F_6 | 1 | 0.05 | F_19 | 1 | 0.06 | T_5 | 1 | 0.10 |
| | 2 | 0.75 | | 2 | 0.02 | | 2 | 0.68 |
| | 3 | 0.09 | | 3 | 0.01 | | 3 | 0.36 |
| | 4 | 0.22 | | 4 | 0.05 | | 4 | 0.71 |
| | 5 | 0.27 | | 5 | 0.06 | | 5 | 0.61 |
| | 6 | 0.44 | | 6 | 0.02 | | 6 | 0.02 |
| | 7 | 0.06 | | 7 | 0.07 | | 7 | 0.89 |
| F_7 | 1 | 0.02 | F_20 | 1 | 0.07 | T_6 | 1 | 0.50 |
| | 2 | 0.61 | | 2 | 0.02 | | 2 | 0.44 |
| | 3 | 0.05 | | 3 | 0.02 | | 3 | 0.27 |
| | 4 | 0.22 | | 4 | 0.06 | | 4 | 0.17 |
| | 5 | 0.23 | | 5 | 0.06 | | 5 | 0.25 |
| | 6 | 0.12 | | 6 | 0.02 | | 6 | 0.93 |
| | 7 | 0.08 | | 7 | 0.07 | | 7 | 0.03 |
| F_8 | 1 | 0.43 | F_21 | 1 | 0.07 | T_7 | 1 | 0.06 |
| | 2 | 0.11 | | 2 | 0.02 | | 2 | 0.35 |
| | 3 | 0.25 | | 3 | 0.01 | | 3 | 0.14 |
| | 4 | 0.34 | | 4 | 0.04 | | 4 | 0.49 |
| | 5 | 0.26 | | 5 | 0.06 | | 5 | 0.37 |
| | 6 | 0.24 | | 6 | 0.04 | | 6 | 0.02 |
| | 7 | 0.30 | | 7 | 0.04 | | 7 | 0.65 |
| F_9 | 1 | 0.15 | F_22 | 1 | 0.43 | T_8 | 1 | 0.58 |
| | 2 | 0.03 | | 2 | 0.04 | | 2 | 0.13 |
| | 3 | 0.04 | | 3 | 0.17 | | 3 | 0.21 |

| | | | | | | |
|---|---|---|---|---|---|---|
| 4 | 0.07 | 4 | 0.20 | 4 | 0.10 |
| 5 | 0.10 | 5 | 0.19 | 5 | 0.18 |
| 6 | 0.04 | 6 | 0.14 | 6 | 0.79 |
| 7 | 0.14 | 7 | 0.30 | 7 | 0.05 |