

Two-Part Factor Mixture Modeling:
Application to an Aggressive Behavior Measurement Instrument

YoungKoung Kim
The College Board
Teachers College, Columbia University

Bengt O. Muthén
University of California, Los Angeles

The research of the first author was sponsored by the Center for Prevention and Early Intervention (P30 MH066247) and jointly funded by the National Institutes of Mental Health and Drug Abuse. The research of the second author was supported by grant K02 AA 00230 from the National Institute on Alcohol Abuse and Alcoholism, by grant 1 R21 AA10948-01A1 from the NIAAA, by NIMH under grant No. MH40859, and by grant P30 MH066247 from the NIDA and the NIMH.

Abstract

This study introduces a two-part factor mixture model as an alternative analysis approach to modeling data where strong floor effects and unobserved population heterogeneity exist in the measured items. As the names suggests, a two-part factor mixture model combines a two-part model, which addresses the problem of strong floor effects by decomposing the data into dichotomous and continuous response components, with a factor mixture model, which explores unobserved heterogeneity in a population by establishing latent classes. Two-part factor mixture modeling can be an important tool for situations where ordinary factor analysis would produce distorted results and can allow researchers to better understand population heterogeneity within groups. Building a two-part factor mixture model involves a consecutive model building strategy that explores latent classes for each part of the data as well as a combination of the two parts. This model building strategy was applied to data from a randomized preventive intervention trial in Baltimore public schools administered by the Johns Hopkins Center for Early Intervention. The proposed model revealed otherwise unobserved subpopulations among the children in the study in terms of their tendency towards and level of aggression. Furthermore, the modeling approach was examined using a Monte Carlo simulation.

Two-Part Factor Mixture Modeling: Application to an Aggressive Behavior Measurement Instrument

Introduction

This paper considers modeling issues arising from the latent variable analysis of items with two common types of complications – data exhibiting strong floor or ceiling effects, which produce highly skewed items, and data arising from several unobserved subpopulations, which produce unobserved heterogeneity. In such situations, conventional factor analysis may give strongly distorted results.

When the first complication – a strong floor effect – is present, a factor analysis measurement model is distorted due to the violation of the multivariate normality assumption and the linearity of the regressions of items on factors. A typical example of strong floor effects is seen in studies of early childhood behavior, with subgroups of children exhibiting high levels of aggressive, hyperactive, impulsive, and inattentive behavior. It is common that items used to measure this type of behavior would show a preponderance of zeros, since the behavior has not yet emerged for many individuals in the population. Two-part modeling of longitudinal data, first introduced by Olsen and Schafer (2001) and applied to intervention studies by Brown, Catalano, Fleming, Haggerty & Abbot (2005), addresses the problem of a preponderance of zeros when analyzing data from abnormal behavior studies. Two-part modeling, as the name suggests, decomposes the distribution of data into two parts – one part which determines whether the response is zero or not and the other part which determines the actual level if non-zero responses occur.

The second complication, unobserved heterogeneity, is often seen in general population samples where both normative and various types of non-normative behavior are exhibited. Factor mixture modeling, which combines factor analysis with a classification of individuals into types in line with latent class analysis, is a useful tool for exploring population heterogeneity (Muthén, 2008; Muthén & Asparouhov, 2006). In longitudinal intervention studies, factor mixture analysis on baseline data can uncover subpopulations that might respond differently to an intervention.

Given the limitations of conventional factor analysis which often cannot handle these two complications properly, this study introduces a two-part factor mixture model as an alternative modeling approach to dealing with data that have strong floor effects for individual items of behavioral measurement and that show heterogeneity. In doing so, the aims are to 1) discuss three model building steps for two-part factor mixture models that combine the components of both two-part and factor mixture models and 2) assess their viability through analyzing the results of the model in a Monte Carlo simulation study. Establishing two-part factor mixture modeling as an important tool for situations where ordinary factor analysis would produce distorted results can reap considerable rewards in practice. Particularly for intervention studies, two-part models hold the potential to allow researchers to better understand population heterogeneity within groups of at-risk children and provide better and more effective intervention techniques that can be tailored to subgroups that exist in a given population.

Method

Two-Part Factor Mixture Model

Two-part factor mixture modeling combines aspects of both factor mixture modeling, which attempts to discover latent classes, and two-part modeling, which has been developed to deal with semi-continuous variables. As an introduction to the methodology, a description of the factor mixture model as well as the two-part model will be provided. This will serve as the background to the later introduction of the combination of these two components into a single two-part factor mixture model.

Factor Mixture Models. Factor mixture models were originally proposed to detect unobserved population heterogeneity (Jedidi, Jagpal, & DeSarbo, 1997; McLachlan & Peel, 2000; McLachlan, Do, & Ambroise, 2004; Mislevy & Verhulst, 1990; Muthén, 2006; Muthén & Asparouhov, 2006; Yamamoto & Gitomer, 1993; Yung 1997). Muthén (2008) provided an overview of these factor mixture models, describing them as hybrid latent variable models. In this overview, latent variable models were categorized into two broad types based on measurement invariance or non-invariance; FMA is classified in the non-invariant measurement branch. Muthén (2008) described how a factor mixture analysis (FMA) model presents a useful generalization for cross-sectional latent variable models, allowing for both the classification of subjects in the form of latent classes and the determination of continuous latent scores within these classes. Compared to the Latent Class Analysis (LCA) that specifies that items are uncorrelated within each latent class, FMA allows the items to have non-zero correlations within each class because the factor in an FMA influences all items. Muthén (2008) points out that an LCA is a special case of FMA where the factor in an FMA is absent and that a variety of FMA models are

possible by including measurement non-invariance in intercept differences and slope differences.

Based on these factor mixture models proposed by Muthén (2008), a factor mixture model for $k= 1, \dots, K$ latent classes can be specified as follows:

$$\begin{aligned} \mathbf{y}_{ik} &= \mathbf{v}_k + \mathbf{\Lambda}_k \boldsymbol{\eta}_{ik} + \boldsymbol{\varepsilon}_{ik} \\ \boldsymbol{\eta}_{ik} &= \boldsymbol{\alpha}_k + \boldsymbol{\zeta}_{ik} \end{aligned} \quad (1)$$

where for class k , \mathbf{y}_{ik} are the individual i 's responses on random variable \mathbf{y} , which is a p vector of observed outcomes; \mathbf{v}_k is a p vector of measurement intercept; $\mathbf{\Lambda}_k$ is a $p \times m$ (number of factor) matrix of factor loadings; $\boldsymbol{\eta}_{ik}$ is an m vector of factor scores and $\boldsymbol{\varepsilon}_{ik}$ is a p vector of residual errors. $\boldsymbol{\alpha}_k$ is an m vector of the intercepts of the factors for each class k .

Factor analysis typically uses the maximum likelihood (ML) estimator. This method, however, can break down when the normality assumption is violated, yielding distorted test statistics and standard errors that can lead to erroneous conclusions (Boomsma & Hoogland, 2001; Muthén & Kaplan, 1985, 1992; Powell and Schafer, 2001; Yuan & Bentler, 1998). Assessing the many robust methods that have been developed to handle such non-normal data, however, Muthén (1989) pointed out that since robust approaches, such as the asymptotically distribution-free (ADF) estimation method, still maintain the linearity assumption, the application of ADF is not appropriate when the linearity of measurement variables are questionable, a situation which often occurs in censored data. In this context, the term censored refers to censoring from below, or left-censoring, which is the same phenomenon observed in data containing a preponderance of zeros. Thus Muthén proposed a Tobit factor analysis approach for the censored data.

Two-Part Models. Two part models are particularly useful in dealing with semi-continuous variables. Similar to left-censored variables, semi-continuous variables have highly skewed distributions with a large portion of observations piled up at a single value – typically zero. According to Olsen and Schafer (2001), a semi-continuous variable is different from one that has been left-censored or truncated, because the zeros are valid self-representing data values, not proxies for negative or missing responses. In practice, semi-continuous variables are frequently found in studies of abnormal behavior, adolescent substance use, and medical expenses.

When a Tobit model is applied to semi-continuous data, Olsen and Schafer (2001) discussed two possible problems. First, when a zero is a valid self-representing data point, the underlying distribution of the censored data does not exist. Thus, the interpretation of parameters – the mean and variance of censored data – can be a problem. Second, in a Tobit model, the censoring mechanism is jointly modeled with the outcome variable generation. When censoring mechanisms and outcome variable generations have separate processes that result in semi-continuous data, the restriction of a Tobit model is not appropriate. As an alternative, Duan et al. (1983) and Manning et al. (1981) used a two-part regression modeling approach to handle data with such a piling up of zeros. Olsen and Schafer (2001) then extended the two-part regression approach to longitudinal settings.

Based on the definition by Olsen and Schafer (2001), a semi-continuous response ranging from zero to $+\infty$, y_{ij} , for an individual $i = 1, \dots, I$ at occasion $j = 1, \dots, J$, can be written as follows:

$$\begin{aligned}
U_{ij} &= \begin{cases} 1 & \text{if } y_{ij} > 0 \\ 0 & \text{if } y_{ij} = 0 \end{cases} \\
V_{ij} &= \begin{cases} g(y_{ij}) & \text{if } y_{ij} > 0 \\ \text{irrelevant} & \text{if } y_{ij} = 0, \end{cases} \tag{2}
\end{aligned}$$

where g is a monotonically increasing function that will make V_{ij} approximately Gaussian. Olsen and Schafer (2001) modeled the semi-continuous responses by using a pair of correlated random-effect models, one for the logit probability of a non-zero response, $U_{ij}=1$ and one for the mean of the continuous responses given that non-zero responses occur, $E(V_{ij} | U_{ij}=1)$. The first part of the model separated “no-use” from “any sort of use” by creating binary indicator variables that reveal any level of use within the previous time. In the second part of the model, continuous indicator variables represent the amount of the usage if “use” occurred. If there was “no-use” on the binary indicator variables, the continuous indicator variable that captures frequency of use was treated as missing. The random coefficients from the two parts were assumed to be jointly normal and possibly correlated.

In the application of the two-part model on abnormal behavior, the dichotomous part at any given time point concerns the engagement in the abnormal activity while the continuous part at any given time point concerns the amount of the activity when the engagement occurred. The engagement in the activity is usually positively correlated with the amount of activity; the higher the probability of engagement is, the higher the expected amount of activity and, conversely, the smaller the probability of engagement is, the lower the expected amount of activity. It is possible, however, for the amount of activity to be high even with a small probability of engagement, but this is a rare occurrence.

Two-Part Factor Mixture Analysis. The two-part modeling approach can be combined with a factor model to deal with a situation where multiple indicators have a preponderance of zeros and the rest of the observations are highly skewed. Thus, the combination of these two ideas takes into account the decomposition of the semi-continuous outcome measurements into a dichotomous response part (i.e., for zero versus non-zero responses) as well as a continuous response part. Figure 1 displays the process of decomposing the skewed distribution of the data into two parts and applying two-part modeling to a factor model which creates a two-part factor model.

Insert Figure 1 about here

Incorporating latent classes in a two-part factor model allows the two-part factor mixture model to explore qualitatively different subpopulations within the data set. The factor mixture model in equation (1) can be decomposed into two parts – one to model the dichotomous response part and the other one to model the continuous response part. First, suppose $(u)_{ik}$ denotes the individual i 's dichotomous response part for class k . The model for the dichotomous part can be written as the following:

$$\begin{aligned}
 y_{ik}^* &= \Lambda_{(u)k} \boldsymbol{\eta}_{(u)ik} + \boldsymbol{\varepsilon}_{(u)ik} \\
 (\mathbf{u})_{ik} &= \begin{cases} 1 & \text{if } y_{ik}^* > 0 \\ 0 & \text{if } y_{ik}^* = 0 \end{cases}, \quad (3)
 \end{aligned}$$

where a p vector of zero or non-zero outcomes $(\mathbf{u})_{ik}$ are observed for each y_{ik}^* which is a set of latent response variables for class k ; $\Lambda_{(u)k}$ is a $p \times m$ matrix of factor loadings for

the dichotomous part for class k ; $\boldsymbol{\eta}_{(u)ik}$ and $\boldsymbol{\varepsilon}_{(u)ik}$ denote an m vector of factor scores and a p vector of residuals, respectively, for the dichotomous part for class k .

Second, similar to the factor mixture model in (1), the model for the individual i 's continuous response part y_{ik} for class k where the observed y_{ik}^* is greater than 0 can be written as follows:

$$y_{ik} = \boldsymbol{v}_k + \boldsymbol{\Lambda}_k \boldsymbol{\eta}_{ik} + \boldsymbol{\varepsilon}_{ik} . \quad (4)$$

Since the continuous response part is usually skewed, the model can use a function g to make y_{ik} normal. Often in practice, logarithm functions are usually employed assuming a log-normal distribution on the continuous response part within each class (Olsen and Schafer, 2001).

The two-part factor mixture model was estimated using a maximum likelihood estimator with robust standard errors in the Mplus program version 4.2 (Muthén & Muthén, 1998-2006). Numerical integration is necessary in maximum-likelihood estimation when a continuous latent variable has categorical indicators as is the case for the dichotomous part of the model.

Three Steps for Building a Two-part Factor Mixture Model

Since two-part factor analysis is a complex undertaking consisting of factor mixture modeling for both the dichotomous outcome part and the continuous outcome part, a stepwise approach should be utilized in constructing a two-part factor mixture model to help avoid model misspecification. This approach involves repeating the same model building strategy three times – once for the dichotomous outcome part by itself, once for the continuous outcome part by itself and finally once for the combination of

the two outcomes. Figure 2 illustrates this multi-step process of building a two-part factor mixture model.

Insert Figure 2 about here

Components of Each Step. Each of the three steps towards building a two-part factor mixture model involves 1) conducting a conventional factor model analysis – either an Exploratory Factor Analysis (EFA) or a Confirmatory Factor Analysis (CFA) to decide the number of factors and 2) specifying a series of models and comparing fit information for each model considered in order to determine the number of latent classes which best captures population heterogeneity.

It is also possible to test models with either class-specific factor loadings or class-specific intercepts/thresholds. During model specification, models with class-specific or class-invariant variances can be compared. For model identification, however, both intercepts and factor means cannot vary across classes simultaneously. Next, the fit of the models with different numbers of classes is compared in order to determine how many classes are needed. Because regularity conditions are not met (McLachlan & Peel, 2000) when comparing mixture models that differ by one class, the traditional chi-square difference test in the form of the likelihood ratio test is not applicable. Instead, information criteria such as the Bayesian Information Criterion (BIC; Schwartz 1978) as well as the bootstrapped Likelihood Ratio Test (BLRT; Nylund, Asparouhov, & Muthén, 2007) have been used as the model selection tool determining the number of classes. Since there are no tests available among likelihood-based tests to compare models that differ in terms of both the number of factors and classes (e.g., a two-factor model with

two classes versus a three-factor model with three classes), the models in this study were compared based on the BIC, the number of parameters and the log likelihood values of the models.

Step One: Dichotomous Component. The first step is to fit a model for only the dichotomous component (i.e., the zero versus non-zero responses) which is Step One in Figure 2. First, an EFA on the dichotomous outcomes is used to study the underlying structure of the data. Based on the results of an EFA, a CFA with a single class follows. Then, more latent classes can be added to the model obtained from the CFA. Given that both classes and factors capture heterogeneity, it should be expected that as the number of classes increases the number of factors needed might decrease. As mentioned earlier, fit indices are used to decide the numbers of classes (e.g., BIC and BLRT).

Step Two: Continuous Component. Step Two in Figure 2 displays the model for the continuous part. In the second step, the same strategy – conducting an EFA, a CFA, and an FMA – can be applied to the continuous outcome component. This is conducted to understand the population heterogeneity of the frequency of use or level of activity in the continuous part of the data.

Step Three: Combination. The final step connects the models found separately in Steps One and Step Two. As displayed in Step Three of Figure 2, this step connects these two modeling components by correlating the factors from the dichotomous component with the factors from the continuous component. Step Three can start with having a single factor from the dichotomous part correlated to a single factor from the continuous part and from there increasing the number of factors from both parts, a process which can serve as an EFA. Then, the latent classes can be added to the two-part

model, and correlated if necessary. Fit indices from all the models under consideration are collected and compared to determine the best-fitting model. The number of factors and latent classes found at Step One and Step Two can provide useful information about the number of factors and classes at Step Three. It should be cautioned, however, that results at Step Three, in terms of the number of factors and classes, might be different from the results of the prior steps because Step Three is the final step connecting the separate two-part analyses done during the prior steps.

In Figure 2, the arrows from the latent class variable to the indicators – i.e. from cu to u and from cy to y – indicate that the item thresholds for the dichotomous part indicators or the intercepts for the continuous part indicators vary across classes. Alternatively, it is possible to allow factor means to vary across classes i.e. allowing arrows from cu to fu and from cy to fy instead of allowing class-specific thresholds and intercepts in the model. For model identification, however, the latent class variables cannot affect both the items and the factors at the same time. In most applications of latent class analysis, the main objective is finding classes that differ with respect to their means or locations (Muthén, 2008). Thus, typical LCA allows the thresholds of categorical indicators or the means of continuous indicators for the latent class variables to vary across classes. Furthermore, a previous study of a factor mixture analysis on tobacco dependence data by Muthén & Asparouhov (2006) found that the approach of allowing latent class measurement parameters – thresholds and intercepts – to vary across classes fitted the data better. On the data for children’s aggressive behavior for the current study, the model with class-specific factor means was also tested and rejected against the model with class-specific thresholds and intercepts. Therefore, this study

chose the model with class-specific item thresholds (for the dichotomous part) and intercepts (for the continuous part).

Data

The data used in this study were obtained from a randomized universal preventive intervention trial on in Baltimore public schools administered by the Johns Hopkins Center for Prevention and Early Intervention.¹ This trial is part of an ongoing research project that the Center has administered since 1985 and has provided the foundation for three generations of school-based preventive intervention field trials and their subsequent follow-ups. In these trials, teacher ratings of each child's aggressive classroom behavior for grades 1-7 were measured. The ratings were made using the Teacher's Observation of Classroom Adaptation Revised (TOCA-R) scaling instrument (Werthamer-Larsson, Kellam & Wheeler, 1991). This rating consists of 10 items, each rated on a six-point scale from 'almost never' to 'almost always'². This study focused on analyzing the pre-intervention data of Cohort 1 – the TOCA-R ratings from the first grade – when the children first entered the intervention trial. Specifically, 527 male students in first grade in 1985 were analyzed. Figure 3 displays the distribution of items that clearly show a strong floor effect. On average, about 50 % of children for each item were in the “almost never” category for each item. Thus, these items cannot be treated as if they were normally distributed.

¹ Formerly the Johns Hopkins Prevention Intervention Research Center (JHU PIRC).

² The ten items are Stubborn, Break Rules, Harms others and Property, Breaks Things, Yells at Others, Take others property, Fights, Lies, Teases Classmates, Trouble Accepting Authority

Insert Figure 3 about here

Results

Using the multi-stage strategy described above, a two-part factor mixture model was fitted to the TOCA-R data. For model estimation, a maximum likelihood estimator with robust standard errors using a numerical integration algorithm was used. To avoid local solutions, a sufficient number of starting random starts was chosen for each step.³ The measurement intercepts (for the continuous part) and thresholds (for the dichotomous part) were allowed to be different across classes in the model. Also for identification purposes, the factor means were fixed at zero in all classes. Since the model estimation was computationally demanding, the factor loadings, variance, and covariance of factors were held equal across classes.

In Step One, the dichotomous response part was examined (i.e., aggression versus non-aggression). [An exploratory factor analysis was conducted and found a 2 factor solution. Based on the EFA result, in Step One, a confirmatory factor model with 2 factors and 1 class was selected for the dichotomous part because the model has the lowest BIC compared to the other models included in Table 1]. Step Two looked at the continuous component of students' level of aggression; it was found that a factor mixture with 1 factor and 2 classes was better than other competing models found in Table 2.

Insert Table 1 about here

³ A minimum of 100 initial stage random sets of starting values and a minimum of 10 final stage optimizations were chosen for each step.

Insert Table 2 about here

Finally in Step Three, the model with 1 factor and 2 classes for each part – the dichotomous component and the continuous component – was selected as the final model based on log likelihood values, BIC and the number of parameters. Table 3 shows that the incorporation of the latent classes improved the model fit compared to the two-part factor model with a single class. Allowing latent classes from the dichotomous and continuous components to covary also seemed to improve the model fit. In the chosen model, therefore, the two latent class variables were allowed to covary. Although in Step One, the model with 2 factors and 1 class was found to be the best solution for the dichotomous part, it was found that in the joint step (Step Three), the model with 2 factors and 1 class for the dichotomous part did not converge due to the correlation between the two factors in the dichotomous part being greater than 1, indicating that the two factors are not statistically distinguishable. This implies that the two-part modeling strategy of integrating both parts can provide different results from modeling only the dichotomous part.

The model where the measurement intercepts and the factor loadings are class-invariant while the factor means are class-specific was also compared to the final model. The measurement invariant model for the continuous part estimated a lower log likelihood (-2776.306 with degree of freedom 56) and a higher BIC (5903.575) than the final model, indicating the measurement invariance for the continuous part did not hold.

Insert Table 3 about here

Interpreting the Results. Table 4 displays the estimated factor loadings of the 10 TOCA-R rating of children's aggressive behavior for the dichotomous part and the continuous part. All factor loadings from the dichotomous part and the continuous part were positive and significant. Thus, the factor found in the dichotomous part can be interpreted as the "propensity to engage in aggressive behavior." The factor from the continuous part can be interpreted as the "propensity to have high aggressive activity levels". The correlation between the two factors was 0.935 and it was significantly different from zero ($p < .01$). Note that the overall size of factor loadings for the dichotomous part is larger than that for the continuous part. In particular, the largest difference in factor loadings between the dichotomous part (1.570) and the continuous part (0.233) was found in item 4 (Break things).

Figure 4 shows the estimated probabilities by latent class for the dichotomous component and the estimated means by latent classes for the continuous component. These plots clearly show the difference in the "propensity to engage in aggressive behavior" as well as the "propensity to have high activity levels" between the groups. The dichotomous component of the model provided the two classes that distinguished two groups of children in terms of their propensity to engage in aggressive behavior – a "High Engagement Group" and a "Low Engagement Group." The endorsements in the "Low Engagement Group" for the dichotomous part were particularly low on item 3 (Harms others and property) and item 4 (Break things), a fact that may suggest that membership in the "Low Engagement Group" versus the "High Engagement Group" seems closely related to physical aggression.

The continuous component of the model that considers the aggressive activity level identified two groups in terms of their propensity to have high activity levels – a “High Activity level Group” versus a “Low Activity Level Group.” In Figure 4, the estimated activity level means for the “High Activity level Group” are higher than the means for the “Low Activity Level Group” across all 10 items. Among 10 items, the mean difference between the “High Activity level Group” and the “Low Activity Level Group” was the largest in item 4 (Break things), item3 (Harms others and property) and item 7 (Fights), indicating physical violence. As seen in the dichotomous component of the model, the items related to physical violence in the continuous component also might play an important role in making class distinctions.

Insert Table 4 about here

Insert Figure 4 about here

Combining the two latent classes from each component of this two-part model, four types of latent class patterns emerged: 1) the Low Engagement/Low Activity Level group, 2) the Low Engagement/High Activity Level group, 3) the High Engagement/Low Activity Level group and 4) the High Engagement/High Activity Level group. Table 5 presents the estimated intercepts and thresholds for the four latent class patterns. The Low Engagement/Low Activity Level group has a lower probability of endorsing each item and lower intercepts for each item compared to the High Engagement/High Activity level group. Although the intercepts of the Low Engagement/High Activity level group

are the same as those in the High Engagement/High Activity level group, the two groups have different probabilities of endorsing the items. Similarly, although the High Engagement/Low Activity Level group has the same intercepts as the Low Engagement/Low Activity Level group, the two groups have different probabilities of endorsing the items.

Insert Table 5 about here

Table 6 displays the counts and proportions of the four patterns. Of the approximately 32% of the students who were in the “High Engagement” group, only 8% fell into the “High Activity Level” group. Thus, about 23% of the students who were in the “High Engagement” group were assigned to the “Low Activity Level” group in the continuous component of the model. Moreover, 68% of the remaining students fell into the “Low Engagement” group. Among these students, 64% show a propensity to have low activity levels by being members of the “Low Activity Level” group. On the other hand, 5% showed a propensity to have high activity levels since they also were members of the “High Activity Level” group. Thus, even with a small probability of engagement the level of aggressive activity can be high although it unlikely to occur.. This indicates that there is a certain group of students who does not show any engagement in aggressive behavior on most of the items but when they do, their level of aggressive activity for those small number of items was very high. Based on the posterior probability of latent class membership, the model classified 13 students in the “Low Engagement/High Activity Level” group. The students in this group received “almost never”, which was the

lowest rating on average about 5 items while they received “almost always” or “very often” for some of items.

Insert Table 6 about here

Figure 5 displays the rating pattern of two students classified into the “Low Engagement/High Activity Level group”. Although Student 1 received “almost never” for 9 items, he received “very often” for item 7 (Fights). Student 2 received almost always for item 1 (Stubborn) and item 9 (Trouble accepting authority) while he received “almost never” for 8 items. Research on types of adolescent aggression has indicated that there is some asymmetry in the high degree of association between the types of aggression (Munoz, Frick, Kimonis & Aucoin, 2008). In other words, two groups of aggressive children are possible when two types of aggression exist: the first group which is highly aggressive and shows both types of aggressive behavior and the other group which is less aggressive overall and shows only one type of aggressive behavior. Thus, the results of the current study suggest that groups based on aggression type, including asymmetric combinations of aggression such as the Low Engagement/High Activity Level group and the High Engagement/Low Activity Level group, can be captured if the characteristics of the data with a preponderance of zeros were taken into account. It will be of great interest to see how the propensities to engage in aggressive behavior and the levels of aggression develop in each group over time.

Insert Figure 5 about here

Comparison to a Regular EFA. It is interesting to note that the BIC from the regular EFA suggested a two-factor solution (BIC for 1 factor solution = 11315.696, BIC for 2 factor solution = 11195.266 and BIC for 3 factor solution = 11206.478). The factor loadings are shown in Table 7. The first factor can be interpreted as a “Verbal Aggression Factor” since the items – Stubborn, Trouble accepting authority, Break rules, Yells at others and Teases classmates – loaded on the first factor. The second factor can be interpreted as a “Physical Aggression Factor” since the following items – Break things, Take others property, Harms others, and Fights – were loaded strongly on the second factor. It was found that the selected two-part factor mixture model with 2 classes fits the data better than a regular two-part factor model i.e. a two-part factor model with a single latent class. Therefore, not only did the two-part factor mixture model capture the common content of observed variables i.e. the aggressiveness factor, it also revealed unobserved population heterogeneity that clustered the children in the study in terms of their propensity to engage in aggressive behavior and their propensity to have aggressive activity levels. By allowing item probabilities to vary across classes, however, the selected two-part factor mixture model found 1 factor and 2 classes indicating that the 2 factors found by EFA should be seen as 1 factor with 2 classes. This shows how conventional factor analysis can give misleading results by ignoring the problem of a preponderance of zeros.

Insert Figure 7 about here

A Monte Carlo Simulation Study

To examine whether the multi-stage strategy used in this application was a reasonable choice or not, a small Monte Carlo simulation study was conducted. Using the Monte Carlo facility in Mplus version 4.2 (Muthén and Muthén, 1998-2006), data with semi-continuous variables were generated and then the simulated data were analyzed through the following three steps to examine how well each step identified the true model: 1) modeling only the dichotomous component, 2) modeling only the continuous component and 3) modeling these two components together.

Data Generation. Similar to the real-data analysis, data with a sample size of 527 and ten semi-continuous variables were considered for the simulation study. As the data generation model of the simulation, a two-part factor mixture model with 1 factor and 2 classes for the dichotomous part and with 1 factor and 2 classes for the continuous part was considered. The data generation model had the same model specification as the final model for the TOCA-R data, i.e. the factor loadings and the factor variances for both the dichotomous and continuous parts were class-invariant. In addition, the factor correlation between the dichotomous and the continuous part was set to 0.9. Factor means for both were set to zero for purposes of model identification. The thresholds for the dichotomous part and the intercepts for the continuous part were specified as class-specific. With respect to the dichotomous part, two latent classes were set to have different item profiles. For latent class 1, the thresholds of the first five items were set to -1.0 and those of the last five items were set to 1.0 while for latent class 2, the thresholds of the first five items and for the last five items were set to 1.0 and -1.0 respectively.

The intercepts for the continuous outcomes in the two latent classes were set as 1, 2, and 3 standard deviations apart. Based on those three sets of intercept differences, therefore, the Mplus Monte Carlo facility generated three types of data – Data 1, Data 2 and Data 3 – composed of ten semi-continuous items with different means. Table 8 both presents the values of the three sets of intercept differences and summarize the three types of data generation. All items were set to have lognormal distributions with standard deviation of 1 but were set to have different means for the three types of the simulated data. 100 replications were conducted for each data type. Approximately 50% of the overall responses to each item were zero. Figure 6 displays the distribution of each item from one of the simulated datasets that clearly shows the preponderance of zeros and the right-skewness with a long tail.

Insert Table 8 about here

Insert Figure 6 about here

Simulation Results. The coverage values for the two-part factor mixture model with 1 factor and 2 classes for the dichotomous part and with 1 factor and 2 classes for the continuous part were found to be reasonable in all three types of data (between 0.79 and 0.99) although the coverage values for some of parameters in Data 1 were relatively small. The bias (between -0.06 and 0.07) and the mean square error (ranging 0.003 and 0.095) for fits of the model to the simulated data were found to be fairly small indicating that the parameters were well recovered. The average class proportions for each data type

were about 0.5 and 0.5 for the two latent classes in both the dichotomous and continuous parts indicating that the class sizes were well recovered since the population values for the class proportions for the two latent classes in both the dichotomous and continuous parts were 0.5 and 0.5 as well.

Table 9 displays the results of the model selection for each step of the multistage model building strategy. The first two steps (Step One and Step Two) towards constructing a two-part factor mixture model entail modeling the dichotomous and continuous parts separately. To determine the number of classes in these two steps, the Bayesian Information Criterion (BIC) and a bootstrapped Likelihood Ratio Test (BLRT) were used. For the dichotomous part of the model, overall BIC performance was good in all three types of data although it was better for both Data 2 and Data 3 than for Data 1. While the lowest values of BIC occurred 100% of the time in the true model for both Data 2 and Data 3, they occurred 78% of the time for Data 1. In addition, the performance of BLRT was better in both Data 2 and Data 3 than in Data 1. Non-significant p-values of BLRT occurred in the true model more than 90% of the time for both Data 2 and Data 3. On the other hand, BLRT selected the true model 76% of the time for Data 1.

For the continuous part, the performance of BIC varied across the type of data indicating that BIC seem sensitive to the intercept difference between two latent classes. For Data 3 where the intercepts of the items in the two latent classes differed by 3 standard deviations, BIC found the true model 100 % of the time. In contrast to Data 3, for Data 1 and Data 2, where the intercepts of the items in two latent classes differed by 1 and 2 standard deviation, BIC failed to identify the true model and selected the model

with 1 factor and 2 classes 100 % of the time. Compared to BIC, BLRT seems less sensitive to the intercept differences between two latent classes. For Data 1, Data 2 and Data 3, BLRT selected the true model in 57%, 56% and 66% of the cases, respectively, indicating that BLRT performance seems consistent across the level of intercept difference between latent classes.

In Step Three, BIC and log likelihood values were evaluated to decide the number of classes since BLRT is not available for a model with more than one latent class variable. In the joint step, the lowest values of BIC occurred 100% of the time in the true model for both Data 2 and Data 3. In contrast, the lowest values of BIC occurred 0% of the time in the correct model for Data 1. The result suggests that the true model can be correctly identified by BIC when the intercepts between latent classes differ by at least 2 standard deviations.

The multi-stage approach to building a two-part factor mixture model indicates that it is possible to get model misspecification if only one part is modeled. The simulation study showed that the intercept differences between the latent classes affected the model selection by BIC. When the intercept difference between latent classes is smaller than 2 standard deviations, BIC performance, especially in Step Two which models the continuous part, might be poor. On the other hand, BLRT seems relatively less sensitive to the intercept differences between latent classes. When the intercept differences between latent classes are smaller than 2 standard deviations, BLRT might provide better information in terms of identifying the true model in Step Two. Thus, not only should caution be taken when attempting to identify the correct model using only a single step approach, the multi-step approach is a way to confirm the correct model if

several competing models are under consideration. In doing so, the multi-step approach also can provide further detail as to how well each model performs under each simulation.

Insert Table 9 about here

Discussion

This paper introduces the two-part factor mixture model as a way to model data that have a preponderance of zeros. The two-part factor mixture model suggests a more flexible framework by allowing the modeling of continuous outcomes and categorical outcomes simultaneously, a strategy not possible with either latent class models or factor models. This modeling approach breaks down a variable into two parts, one that identifies if a behavior is observed or not, and another part that describes the extent to which the behavior exists. Incorporating a two-part modeling approach into the factor mixture model also significantly reduced the computational difficulty. This paper includes an application of this model exploring aggressive behavior in first grade children. Results showed that there are four patterns: the Low Engagement/Low Activity Level group, the Low Engagement/High Activity Level group, the High Engagement/Low Activity Level group and the High Engagement/High Activity Level group. Since these models involve several parts, a model building strategy was suggested as a way to specify a two-part factor mixture model and replicated through a Monte Carlo simulation.

Clearly, there are more avenues for further study beyond the scope of this paper. Covariates can be included into the model, such as family background, gender, and race. The advantage of the two-part model is its ability to examine how differently these covariates affect the dichotomous part and the continuous part. In this study, since only a sample of first grade students in the fall were analyzed, the sizes of some latent classes were relatively small. Extending the sample, for example to include spring first grade data, might help the model find unobserved subpopulations.

Moreover, although the multi-stage strategy found to be reasonable in building a two-part model, there still are modeling questions that need to be investigated using simulation studies. First, given that this study used a limited type of a two-part factor mixture model, various types of two-part factor mixture models such as a model with class-varying factor means versus model with class-varying intercepts of the items should be compared in order to evaluate the performance of the model. Furthermore, these two-part factor mixture models should be compared to alternative and possibly less complicated latent class models as well. Second, when the log normality assumption, which was assumed in the continuous response part of the two-part factor mixture model, is violated, it is possible that class enumeration based on BIC can fail (Nylund et al, 2007). Therefore, a simulation study could examine the effect of the violation of the within class log normality assumption and how it affects class enumeration. Third, although the current simulation study showed that the BIC would fail to identify the correct number of classes in the joint step when the intercept difference between two latent classes is not large, further simulation studies should be conducted to examine the performance of BIC at the joint step in more detail by looking at various conditions that

affect class enumeration such as sample size variation, number of items, different model settings, etc. Fourth, the effect of model constraints such as class-invariant factor variances that the current study employed on the class enumeration should be investigated by a simulation comparing other models without the constraints.

Ultimately, the furtherance of this model can help better understand how to treat at-risk children before their abnormal behavior becomes manifest. Especially since signs of abnormal behavior in children can predict serious developmental problems that can afflict the later ability of children to adapt and adjust to society as adults. Studies have found that high levels of disruptive behavior during childhood are associated with negative outcomes in adolescence and adulthood (such as risk for school dropouts, academic difficulties, juvenile delinquency, etc.). Thus, it is important that a methodology that accurately predicts the outcome of behavioral therapies on at-risk children be developed. This would allow for a more effective means of conducting the many widely employed interventions aimed at reducing the level disruptive behavior such as peer and teacher-mediated behavioral interaction. Knowledge of the eventual trajectory of adult abnormal behavior can provide clearer insight into the proper treatment of an individual child at any point along the child's development path. Especially critical are the first stages of development since interventions can be adjusted or fine-tuned early-on to match the specific trajectory of the child. As an advancement upon existing modeling strategies, the proposed two-part factor mixture model can uncover these children that would have gone unnoticed. Based on the results of this study, it will be of great interest to further study a Latent Transition two-part factor mixture analysis analyzing all of the time points beyond the baseline time point. This may help provide

the most effective intervention methods for the children and can ultimately lead to a more successful intervention.

References

- Boomsma, A., & Hoogland, J. J. (2001). The robustness of LISREL modeling revisited. In R. Cudeck, S. du Toit, & D. Sörbom, (Eds.), *Structural Equation Modeling: Present and future* (pp. 139-168). Lincolnwood, IL: Scientific Software International.
- Brown, E.C., Catalano, C.B., Fleming, C.B., Haggerty, K.P. & Abbot, R.D. (2005). Adolescent substance use outcomes in the Raising Healthy Children Project: A two-part latent growth curve analysis. *Journal of Consulting and Clinical Psychology, 73*, 699-710
- Duan, N., Manning, W. G., Morris, C. N., & New house, J. P. (1983). A comparison of alternative models for the demand for medical care. *Journal of Business and Economic Statistics, 1*, 115-126.
- Jedidi, K., Jagpal, H. S., & DeSarbo, W. S. (1997). STEMM: A general finite mixture structural equation model. *Journal of Classification, 14*, 23-50.
- McLachlan, G. J., & Peel, D. (2000). *Finite mixture models*. New York: Wiley & Sons.
- McLachlan, G. J., Do, K.A., & Ambrose, C. (2004). *Analyzing microarray gene expression data*. New York: Wiley & Sons.
- Mislevy, R. J., & Verhulst, N. (1990). Modeling item responses when different subjects employ different solution strategies. *Psychometrika, 55*, 195-215.
- Munoz, L. C., Frick, P. J., Kimonis E. R. & Aucoin, K. J. (2008). Types of aggression, responsiveness to provocation, and callous-unemotional traits in detained adolescents. *Journal of Abnormal Child Psychology, 36*, 15-28.

- Muthén, B. O. (1989). Tobit factor analysis. *British Journal of Mathematical and Statistical Psychology*, *42*, 241-250.
- Muthén, B. O. (2006). Should substance use disorders be considered as categorical or dimensional? *Addiction*, *101*, 6-13.
- Muthén, B. (2008). Latent variable hybrids: Overview of old and new models. In Hancock, G. R., & Samuelsen, K. M. (Eds.), *Advances in latent variable mixture models*, 1-24. Charlotte, NC: Information Age Publishing, Inc.
- Muthén, B. O., & Asparouhov, T. (2006). Item response mixture modeling: Application to tobacco dependence criteria. *Addictive Behaviors*, *31*, 1050-1066.
- Muthén, B. O., & Kaplan, D. (1985). A comparison of some methodologies for the factor analysis of non-normal Likert variables. *British Journal of Mathematical and Statistical Psychology*, *38*, 171-189.
- Muthén, B. O., & Kaplan, D. (1992). A comparison of some methodologies for the factor analysis of non-normal Likert variables: A note on the size of the model. *British Journal of Mathematical and Statistical Psychology*, *45*, 19-30.
- Muthén, L., & Muthén, B. O. (1998-2006). *Mplus User's Guide*. Fourth Edition. Los Angeles, CA: Muthén & Muthén.
- Nylund, K. L., Asparouhov, T., & Muthén, B. O. (2007). Deciding on the number of classes in latent class analysis and growth mixture modeling: A Monte Carlo simulation study. *Structural Equation Modeling: A Multidisciplinary Journal*, *14*, 535-569.

- Olsen, M. K., & Schafer, J. L. (2001). A two-part random effects model for semicontinuous longitudinal data. *Journal of the American Statistical Association*, *96*, 730-745.
- Powell, D. A., & Schafer, W. D. (2001). The robustness of the likelihood ratio chi-square test for structural equation models: A meta-analysis. *Journal of Educational and Behavioral Statistics*, *26*, 105-132.
- Schwartz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, *6*, 461-464
- Werthamer-Larsson, L., Kellam S., & Wheeler, L. (1991). Effect of first-grade classroom environment on shy behavior, aggressive behavior and concentration problems. *American Journal of Community Psychology*, *19*, 585-602.
- Yamamoto, K., & Gitomer, D. H. (1993). Application of a HYBRID model to a test of cognitive skill representation. In N. Fredriksen, R. Mislevy, & I. Bejar (Eds.), *Test theory for a new generation of tests*. Hillsdale, NJ: Erlbaum.
- Yuan, K.-H., & Bentler, P. M. (1998). Normal theory based test statistics in structural equation modeling. *British Journal of Mathematical and Statistical Psychology*, *51*, 289-309
- Yung, Y. F. (1997). Finite mixtures in confirmatory factor analysis models. *Psychometrika*, *62*, 297-330.

Table 1 Result of FMA: step 1 – dichotomous part

Model	Log likelihood	#parameters	BIC
CFA_1f1c	-2369.794	20	4864.932
CFA_2f1c	-2357.190	23	4858.525
FMA_1f2c with no class-specific variance	-2343.326	31	4880.935
FMA_2f2c with no class-specific variance	-2333.956	34	4880.996

Table 2 Result of FMA: step2 – continuous part

Model	Log likelihood	#parameters	BIC
CFA_1f1c	-488.740	30	1160.152
CFA_2f1c	-462.092	34	1131.211
FMA_1f2c with no class-specific variance	-397.959	42	1051.657
FMA_2f2c with no class-specific variance	-393.512	45	1061.031

Table 3 Result of FMA: step3 – joint analysis of dichotomous and continuous parts

Model	Log likelihood	#parameters	BIC
Two-part Factor model	-2769.178	51	5857.983
FMA_Y1f2c_U2f1c	nc	nc	nc
FMA_Y1f2c_U1f2c	-2669.339	73	5796.183
FMA_Y1f2c_U1f2c covarying cy and cu	-2665.564	74	5794.900

Note:

cy: latent class from the continuous part

cu: latent class from the dichotomous part

nc: nonconverged

Table 4 Factor loadings and standard errors for the dichotomous part and the continuous part

Dichotomous part	Estimates	S.E.	Est./S.E.
Stubborn	1.000	0.000	0.000
Break rules	1.598	0.261	6.124
Harms others and property	1.765	0.315	5.609
Break things	1.570	0.471	3.335
Yells at others	1.238	0.176	7.025
Take others property	1.387	0.224	6.187
Fights	1.308	0.246	5.323
Lies	1.389	0.274	5.076
Tease classmates	1.168	0.177	6.599
Trouble accepting authority	1.523	0.298	5.104

Continuous part	Estimates	S.E.	Est./S.E.
Stubborn	1.000	0.000	0.000
Break rules	1.070	0.112	9.517
Harms others and property	0.808	0.126	6.395
Break things	0.233	0.065	3.572
Yells at others	1.035	0.121	8.526
Take others property	0.669	0.109	6.164
Fights	0.843	0.151	5.577
Lies	0.724	0.139	5.230
Tease classmates	0.819	0.127	6.437
Trouble accepting authority	1.016	0.101	10.055

Table 5 Estimated intercepts and thresholds of four latent class patterns

Intercepts	Low Engagement/ Low Activity Level			Low Engagement/ High Activity Level		
	Estimate	S.E.	Est./S.E.	Estimate	S.E.	Est./S.E.
Item 1	0.967	0.021	45.444	1.281	0.073	17.463
Item 2	0.922	0.022	42.542	1.416	0.062	22.977
Item 3	0.756	0.019	39.056	1.283	0.076	16.84
Item 4	0.740	0.014	53.404	1.488	0.054	27.727
Item 5	0.818	0.021	38.938	1.283	0.069	18.656
Item 6	0.818	0.023	36.052	1.353	0.066	20.627
Item 7	0.770	0.022	35.587	1.327	0.074	18.013
Item 8	0.836	0.024	34.263	1.182	0.101	11.719
Item 9	0.909	0.022	42.238	1.233	0.068	18.146
Item 10	0.841	0.026	32.688	1.242	0.081	15.289
Thresholds						
Item 1	-0.542	0.205	-2.649	-0.542	0.205	-2.649
Item 2	-0.999	0.384	-2.598	-0.999	0.384	-2.598
Item 3	2.252	0.636	3.538	2.252	0.636	3.538
Item 4	3.621	0.882	4.105	3.621	0.882	4.105
Item 5	0.673	0.414	1.627	0.673	0.414	1.627
Item 6	1.658	0.668	2.481	1.658	0.668	2.481
Item 7	0.893	0.295	3.026	0.893	0.295	3.026
Item 8	0.758	0.322	2.358	0.758	0.322	2.358
Item 9	-0.219	0.350	-0.624	-0.219	0.350	-0.624
Item 10	0.906	0.250	3.620	0.906	0.250	3.620

Table 5 (continued) Estimated intercepts and thresholds of four latent class patterns

Intercepts	High Engagement/ Low Activity Level			High Engagement/ High Activity Level		
	Estimate	S.E.	Est./S.E.	Estimate	S.E.	Est./S.E.
Item 1	0.967	0.021	45.444	1.281	0.073	17.463
Item 2	0.922	0.022	42.542	1.416	0.062	22.977
Item 3	0.756	0.019	39.056	1.283	0.076	16.840
Item 4	0.740	0.014	53.404	1.488	0.054	27.727
Item 5	0.818	0.021	38.938	1.283	0.069	18.656
Item 6	0.818	0.023	36.052	1.353	0.066	20.627
Item 7	0.770	0.022	35.587	1.327	0.074	18.013
Item 8	0.836	0.024	34.263	1.182	0.101	11.719
Item 9	0.909	0.022	42.238	1.233	0.068	18.146
Item 10	0.841	0.026	32.688	1.242	0.081	15.289
Thresholds						
Item 1	-2.421	0.636	-3.806	-2.421	0.636	-3.806
Item 2	-3.569	0.607	-5.884	-3.569	0.607	-5.884
Item 3	-3.209	1.152	-2.785	-3.209	1.152	-2.785
Item 4	-2.729	2.175	-1.254	-2.729	2.175	-1.254
Item 5	-2.692	0.693	-3.886	-2.692	0.693	-3.886
Item 6	-3.248	0.79	-4.109	-3.248	0.790	-4.109
Item 7	-1.063	0.327	-3.248	-1.063	0.327	-3.248
Item 8	-1.416	0.333	-4.246	-1.416	0.333	-4.246
Item 9	-2.573	0.613	-4.196	-2.573	0.613	-4.196
Item 10	-1.731	0.844	-2.052	-1.731	0.844	-2.052

Table 6 Class counts and proportions based for the latent class patterns based on the estimated model

Latent Class from Dichotomous component	Latent Class from Continuous component	Counts	Proportions
Low Engagement	Low Activity Level	336.745	0.639
Low Engagement	High Activity Level	25.131	0.048
High Engagement	Low Activity Level	121.737	0.231
High Engagement	High Activity Level	43.387	0.082

Table 7 Factor loadings from EFA 2-factor solution

Item	Factor1	Factor2
Stubborn	0.909	-0.053
Break rules	0.605	0.332
Harms others and property	0.212	0.767
Break things	0.014	0.900
Yells at others	0.530	0.409
Take others property	0.165	0.772
Fights	0.381	0.551
Lies	0.450	0.446
Tease classmates	0.493	0.377
Trouble accepting authority	0.787	0.129

Table 8 Summary of data generation for simulation

	Data 1	Data 2	Data 3
Intercepts difference between two classes	1 SD difference Class 1: -2.5 Class 2: -1.5	2 SD difference Class 1: -2.5 Class 2: -0.5	3 SD difference Class 1: -2.5 Class 2: 0.5
Distribution of items	Lognormal ($\mu=-2, \sigma=1$)	Lognormal ($\mu=-1.5, \sigma=1$)	Lognormal ($\mu=-1, \sigma=1$)
Model	Two part factor mixture model – 1 factor 2 classes for both dichotomous part and continuous part		
Sample size	527	527	527
Number of item	10	10	10
Number of replications	100	100	100

Table 9 Model selection by BIC and BLRT : Percentage of times the lowest value of BIC and Percent of times of a non-significant p -value selected for BLRT (Bolded rows represent the true k-class model for the given model)

	Data1		Data2		Data3	
	BIC	BLRT	BIC	BLRT	BIC	BLRT
Step1: Dichotomous part						
1fu_1c	0	0	0	0	0	0
1fu_2c	78	76	100	92	100	94
1fu_3c	22	24	0	8	0	6
Step2: Continuous part						
1fy_1c	100	26	100	5	0	0
1fy_2c	0	57	0	56	100	66
1fy_3c	0	17	0	39	0	34
Step3: Joint Part						
1fu_1c – 1fy_1c	0	-	0	-	0	-
1fu_2c – 1fy_2c	0	-	100	-	100	-
1fu_3c – 1fy_3c	100	-	0	-	0	-

Figure 1 Process of decomposing the skewed data into two parts and applying two-part modeling to a factor model

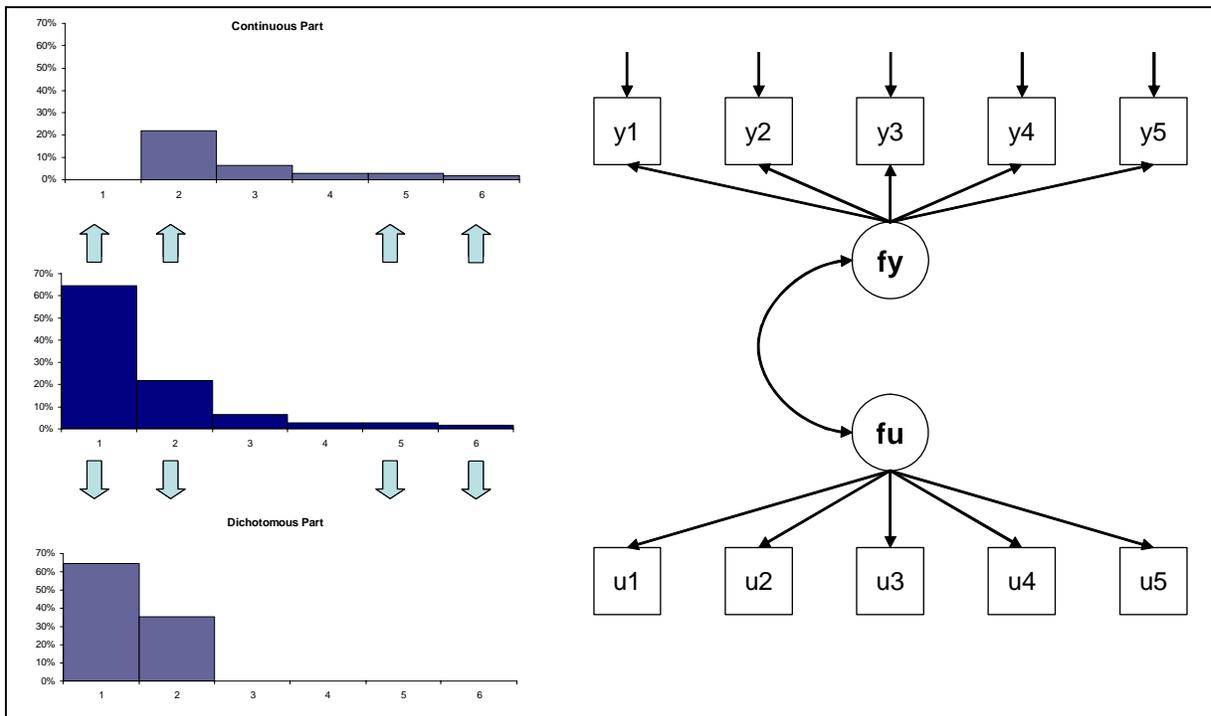


Figure 2 Illustration of multi-stage strategy to build a two-part factor mixture model

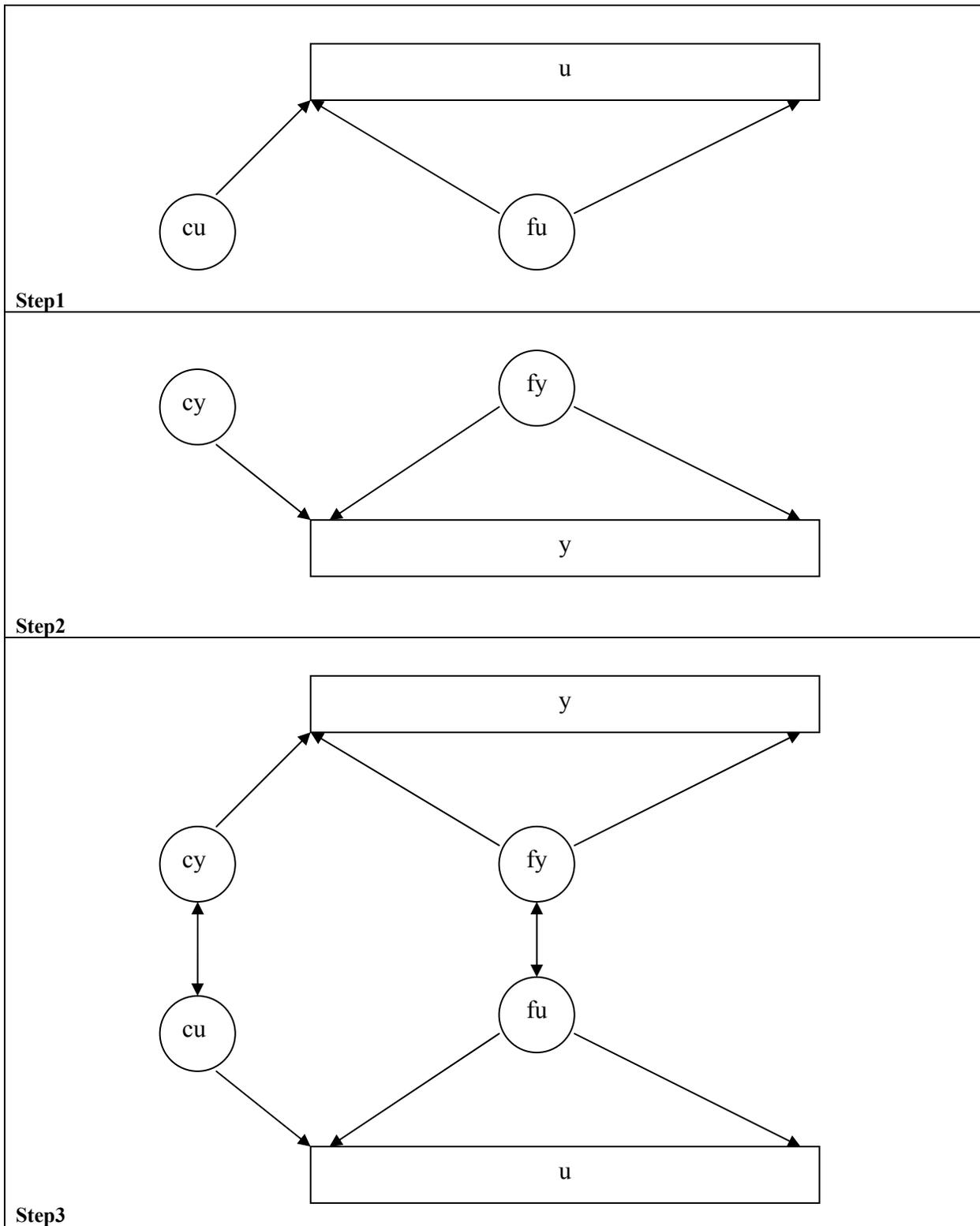


Figure 3 Distribution of 10 TOCA-R items

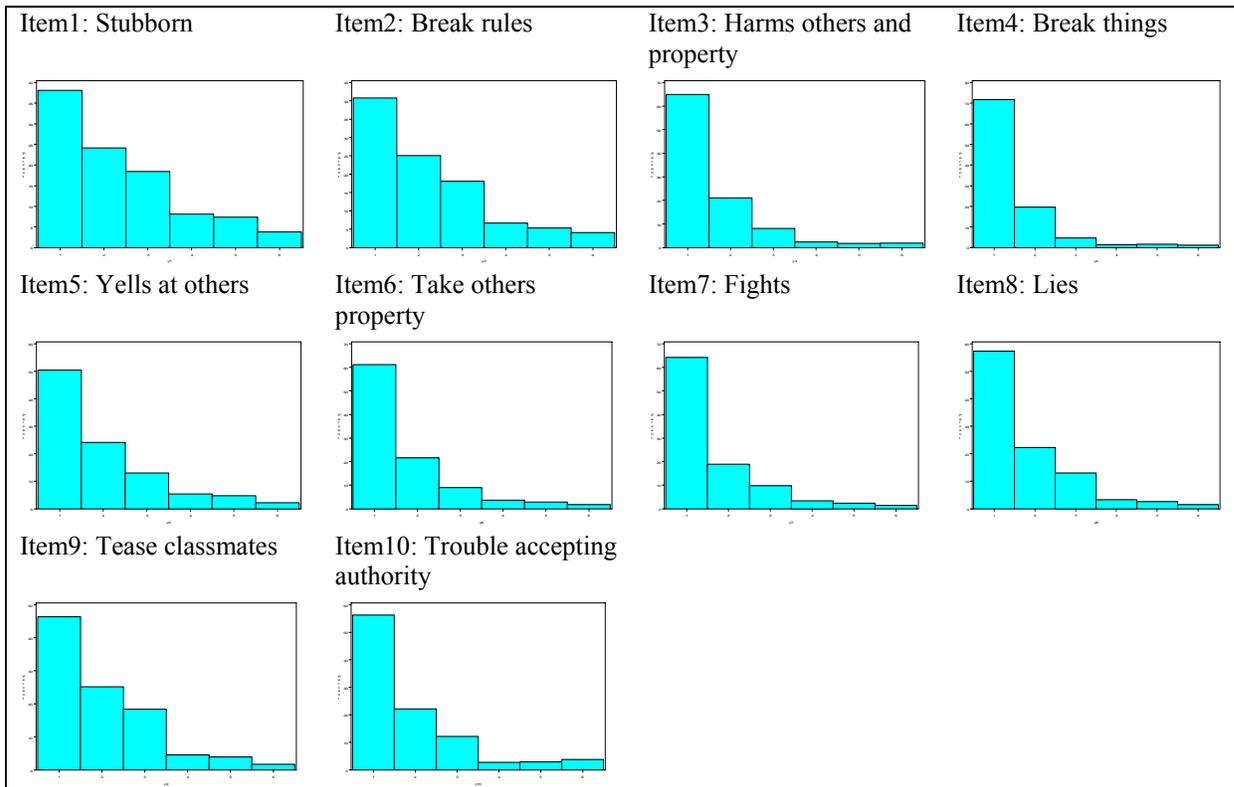


Figure 4 Estimated probabilities for the dichotomous part and the estimated means for the continuous part

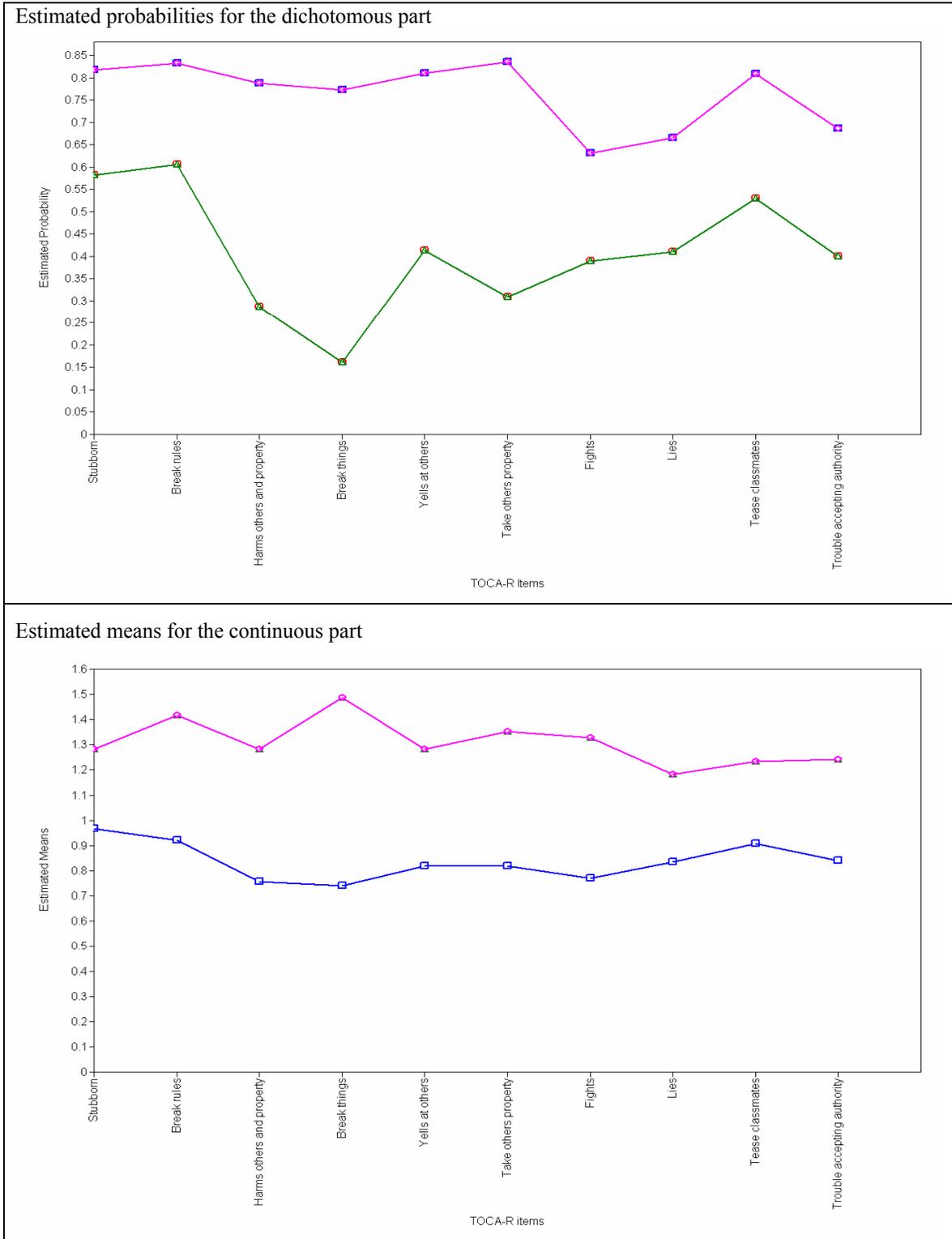
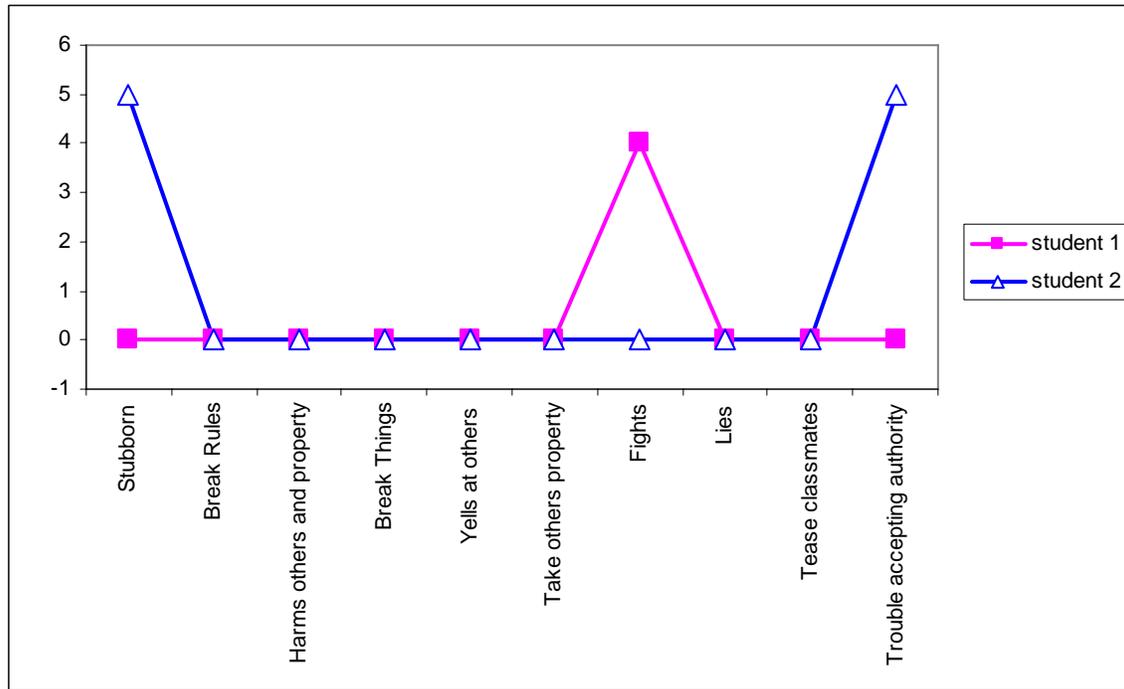


Figure 5 TOCA-R items for two students in the Low Engagement/High Activity Level group



Note: Y-axis 0 = almost never, 1 = rarely, 2 = sometimes, 3 = often, 4 = very often, 5 = almost always

Figure 6 Distribution of 10 items from simulated data

