

Volume 39, Issue 2

Identifying finite mixture models in the presence of moment-generating function: application in medical care using a zero-inflated binomial model

Hiroaki Masuhara

Faculty of Economics and Law, Shinshu University

Abstract

This study presents a simple method to identify the parameters in finite mixture models when a moment-generating function (MGF) is present. We obtain the model conditions using a zero-inflated binomial model, a simple form of the finite mixture binary model, and analyze the results using the Monte Carlo simulation. Using the zero-inflated and standard binomial models, we compare the marginal effects of health care usage.

This work was supported by Shinshu University Grant to advanced sciences. I am grateful to the Editor and three referees for insightful comments and suggestions for improvement. The authors would like to thank Enago (www.enago.jp) for the English language review.

Citation: Hiroaki Masuhara, (2019) "Identifying finite mixture models in the presence of moment-generating function: application in medical care using a zero-inflated binomial model", *Economics Bulletin*, Volume 39, Issue 2, pages 1529-1537

Contact: Hiroaki Masuhara - masuhara@shinshu-u.ac.jp.

Submitted: April 03, 2018. **Published:** June 15, 2019.

1. Introduction

Finite mixture models are widely used in applied econometrics as they are semi-parametric and flexible (Deb and Pravin, 1997; Deb and Trivedi, 2002; Winkelmann, 2004). These models assume the sample of individuals comes from a population containing a finite number of latent classes and that each element is drawn from one of these latent subpopulations or strata. Let y be a random variable and $f(y)$ be its probability density function. Then, a finite mixture model is such that: $f(y | \mathbf{x}, \Theta) = \sum_{j=1}^J p^{(j)} f^{(j)}(y | \mathbf{x}, \theta^{(j)})$, where $p^{(j)}$ is a proportion of j th component ($\sum_{j=1}^J p^{(j)} = 1$), $f^{(j)}(y)$ is a density of j th component, \mathbf{x} is a vector of regressors, $\theta^{(j)}$ is a vector of parameters of j th component, and $\Theta \equiv (\theta^{(1)}, \dots, \theta^{(J)}, p^{(1)}, \dots, p^{(J-1)})'$. This means that the finite mixture model analyzes J types of individuals.

However, some finite mixture models—such as finite mixture *cross-sectional* binomial (probit or logit) models—are not estimated because their parameters are not identified. Teicher (1960) and Blischke (1964) indicated this research gap and presented a sufficient condition for the identifiability of parameters. Their results summarized that the J component mixture with T Bernoulli trials is identifiable if $J \leq (T + 1)/2$. In other words, a two-component finite mixture *cross-sectional* probit or logit model is not identifiable whereas a *panel* probit or logit model with $T \geq 3$ is identifiable. Kasahara and Shimotsu (2014) demonstrated that in finite mixture binomial models the number of components can be non-parametrically identified if $T \geq 2$, and the mixing proportions and distribution of components can be identified when $T \geq 3$.

However, previous studies have not confirmed whether the finite mixture model can be identified using data other than the panel binary data. This paper presents a simple method to identify finite mixture models given the presence of a moment-generating function (MGF). This study demonstrates that a finite mixture model is identifiable if the Jacobian determinant of the joint moments and sample moments is not zero. Although we assume the MGF, *finite* discrete distributions (such as binomial or multinomial variables) always have the MGF. Therefore, our method will be useful in many applied econometric fields.

This study is organized as follows; Section 2 proposes the method to identify a finite mixture model and considers the identifiability of various discrete distributions. Moreover, we propose a zero-inflated binomial model as the simplest finite mixture binomial model. Section 3 depicts the results of the Monte Carlo simulation of the zero-inflated binomial model. Section 4 applies the zero-inflated binomial model using health care data. Section 5 concludes the paper.

2. Discrete multivariate finite mixture models

2.1 Identifying finite mixture models

Let $\mathbf{y} \equiv (y_1, y_2, \dots, y_K)'$ be a $K \times 1$ vector of random variables. To simplify the calculation, we consider the case of *discrete* random variables. If y_k is *finite*, an MGF always exists. The MGF, thus, takes the following form:

$$\begin{aligned} M_{y_1, y_2, \dots, y_K}(t_1, t_2, \dots, t_K) \\ = \sum_{y_1}^{\bar{y}_1} \sum_{y_2}^{\bar{y}_2} \dots \sum_{y_K}^{\bar{y}_K} e^{t_1 y_1 + t_2 y_2 + \dots + t_K y_K} f(y_1, y_2, \dots, y_K | \Theta), \end{aligned} \quad (1)$$

where $f(y_1, y_2, \dots, y_K | \Theta)$ is a probability mass function (PMF), Θ is a parameter to be estimated, and \bar{y}_k is an upper bound of y_k . To simplify discussion, we omit the regressors \mathbf{x} . Moreover, $\mu'_{r_1, r_2, \dots, r_K}$ is the (r_1, r_2, \dots, r_K) th joint moment around the origin of coordinates, which gives

$$\mu'_{r_1, r_2, \dots, r_K} \equiv \sum_{y_1}^{\bar{y}_1} \sum_{y_2}^{\bar{y}_2} \dots \sum_{y_K}^{\bar{y}_K} y_1^{r_1} y_2^{r_2} \dots y_K^{r_K} f(y_1, y_2, \dots, y_K | \Theta). \quad (2)$$

Then, we obtain the following relation:

$$\mu'_{r_1, r_2, \dots, r_K} = \frac{\partial^{(r_1 + r_2 + \dots + r_K)}}{\partial t_1^{r_1} \partial t_2^{r_2} \dots \partial t_K^{r_K}} M_{y_1, y_2, \dots, y_K}(t_1, t_2, \dots, t_K) \Big|_{t_1 = t_2 = \dots = t_K = 0}. \quad (3)$$

In a K -variate finite mixture model, the PMF can be written as

$$f(y_1, y_2, \dots, y_K | \Theta) = \sum_{j=1}^J p^{(j)} f^{(j)}(y_1, y_2, \dots, y_K | \theta^{(j)}), \quad (4)$$

where $f^{(j)}(\cdot)$ is a component PMF, J is the number of components, $p^{(j)} \in (0, 1)$ is a proportion of j th component that satisfies $\sum_{j=1}^J p^{(j)} = 1$, $\Theta \equiv (\theta^{(1)}, \dots, \theta^{(J)}, p^{(1)}, \dots, p^{(J-1)})'$, and $\theta^{(j)} \equiv (\theta_1^{(j)}, \theta_2^{(j)}, \dots, \theta_M^{(j)})'$, where M is the number of parameters of each component. When including regressors \mathbf{x} , the parameter $\theta^{(j)}$ conditional on \mathbf{x} is regarded as a link function $F^{(j)}(\mathbf{x}'\boldsymbol{\beta}^{(j)})$ such as logit or probit distribution, where $\boldsymbol{\beta}^{(j)}$ is a parameter vector of regressors \mathbf{x} of j th component. Therefore, when $\boldsymbol{\beta}$ is identifiable with respect to the link function $F(\cdot)$ and $\theta^{(j)}$ is identifiable within the FM model, $\boldsymbol{\beta}^{(j)}$ is also identifiable. Thus, we only consider the case without regressors for simplicity. Since there exist *sample* joint moments, such as $E[y_1^{r_1} y_2^{r_2} \dots y_K^{r_K}]$, we obtain the following proposition:

Proposition 1. *When the Jacobian determinant of the simultaneous equations composed of joint moments obtained by (2) is not zero, the K -variate and J -component finite mixture model is identifiable.*

Proof. The simultaneous equations of a J -component finite mixture model take the following form:

$$\begin{aligned}
& \sum_{j=1}^J p^{(j)} \mu_{1,0,\dots,0}^{(j)'}(\boldsymbol{\theta}_j) = \text{E}[y_1], \\
& \quad \vdots \\
& \sum_{j=1}^J p^{(j)} \mu_{0,0,\dots,1}^{(j)'}(\boldsymbol{\theta}_j) = \text{E}[y_K], \\
& \sum_{j=1}^J p^{(j)} \mu_{2,0,\dots,0}^{(j)'}(\boldsymbol{\theta}_j) = \text{E}[y_1^2], \\
& \quad \vdots \\
& \sum_{j=1}^J p^{(j)} \mu_{1,0,\dots,1}^{(j)'}(\boldsymbol{\theta}_j) = \text{E}[y_1 y_K], \\
& \quad \vdots \\
& \sum_{j=1}^J p^{(j)} \mu_{r_1,r_2,\dots,r_K}^{(j)'}(\boldsymbol{\theta}_j) = \text{E}[y_1^{r_1} y_2^{r_2} \cdots y_K^{r_K}]. \tag{5}
\end{aligned}$$

These simultaneous equations contain a $(M \times J + J - 1) \times 1$ vector of parameters $\boldsymbol{\Theta}$. Since we assume that the Jacobian determinant of the joint moments and sample moments is not zero and the rank of the Jacobian is $(M \times J + J - 1)$, the K -variate and J -component finite mixture model is identifiable. \square

When we use *finite* discrete variables, this proposition has a great advantage. Since all statistical models using these variables have an MGF, we can confirm the number of identified components of finite mixture models with this proposition.

2.2 Examples

Let us consider a two-component ($J = 2$) finite mixture tri-variate binomial model ($K = 3$). Its component PMF takes the following form:

$$f^{(j)}(y_1, y_2, y_3) = \prod_{k=1}^3 \left[1 - F_k^{(j)}\right]^{1-y_k} \left[F_k^{(j)}\right]^{y_k}. \tag{6}$$

In this case, the number of parameters to be estimated is $3 \times 2 + 1 = 7$. This model can, thus, be identified because the Jacobian determinant of the MGF and sample moments is $(p^{(1)})^3 (1 - p^{(1)})^3 \prod_{k=1}^3 \left[F_k^{(1)} - F_k^{(2)}\right]^2$ and is not zero

when $F_k^{(1)} \neq F_k^{(2)}$. This result is essentially the same as the results of Teicher (1963), Blischke (1964), and van Wieringen (2005).

Moreover, from Proposition 1, we also investigate some “unrealistic” finite discrete distribution. For example, consider the three-component ($J = 3$) bivariate binomial distributions ($K = 2$, y_1 and y_2). Its Bernoulli trials of y_1 is three and that of y_2 is two. Then, its component PMF obtains

$$f^{(j)}(y_1, y_2) = \binom{3}{y_1} [1 - F_1^{(j)}]^{3-y_1} [F_1^{(j)}]^{y_1} \times \binom{2}{y_2} [1 - F_2^{(j)}]^{2-y_2} [F_2^{(j)}]^{y_2}. \quad (7)$$

One of the determinants of the simultaneous equations of joint moments and sample moments is

$$\begin{aligned} & 186,624 \left(F_1^{(1)} - F_1^{(2)}\right) \left(F_1^{(2)} - F_1^{(3)}\right) \left(F_1^{(3)} - F_1^{(1)}\right) \\ & \times \left[F_1^{(1)} \left(F_2^{(2)} - F_2^{(3)}\right) + F_1^{(2)} \left(F_2^{(3)} - F_2^{(1)}\right) + F_1^{(3)} \left(F_2^{(1)} - F_2^{(2)}\right) \right]^4 \\ & \times \left(p^{(1)}\right)^2 \left(p^{(2)}\right)^2 \left(1 - p^{(1)} - p^{(2)}\right)^2. \end{aligned} \quad (8)$$

Therefore, if not $F_1^{(j)} = F_1^{(k)}$ ($j \neq k$), $F_2^{(1)} = F_2^{(2)} = F_2^{(3)}$, $p^{(j)} = 0$ ($j = 1, 2$), or $p^{(1)} + p^{(2)} = 1$, this model with three components is identifiable.

2.3 Zero-inflated binomial models

A simplified finite mixture model is a zero-inflated binomial model. This model treats zero-valued observations as special and its PMF takes the following form:

$$f(y_1, y_2, \dots, y_K) = p^{(1)} \mathbf{I}(y_k = 0, \forall k) + \sum_{j=2}^J p^{(j)} f^{(j)}(y_1, y_2, \dots, y_K). \quad (9)$$

From (9), we can observe that the first component PMF is concentrated at zero with a probability of one; therefore, the zero observations are inflated. When the component PMF follows a Bernoulli distribution, the equation becomes the following:

$$f(y_1, y_2, \dots, y_K) = p^{(1)} \mathbf{I}(y_k = 0, \forall k) + \sum_{j=2}^J p^{(j)} \prod_{k=1}^K [1 - F_k^{(j)}]^{1-y_k} [F_k^{(j)}]^{y_k}. \quad (10)$$

In this zero-inflated binomial model, we obtain the following proposition:

Proposition 2. *A bivariate zero-inflated binomial model is identifiable.*

Proof. Without loss of generality, a bivariate zero-inflated binomial model means $K = 2$ and $J = 2$ in (10), and its PMF is the following:

$$f(y_1, y_2) = p^{(1)}I(y_1 = y_2 = 0) + (1 - p^{(1)}) \prod_{k=1}^2 [1 - F_k]^{1-y_k} [F_k]^{y_k}, \quad (11)$$

where $F_k \equiv F_k^{(2)}$ in (10). The relation between the joint moments and sample moments is the following:

$$(1 - p^{(1)}) F_1 = E[y_1], \quad (1 - p^{(1)}) F_2 = E[y_2], \quad (1 - p^{(1)}) F_1 F_2 = E[y_1 y_2].$$

The Jacobian determinant of simultaneous equations is $(1 - p^{(1)})^2 F_1 F_2$. Therefore, a bivariate zero-inflated binomial model is identifiable if $0 < p^{(1)}, F_1, F_2 < 1$. \square

In a zero-inflated binomial model, a univariate mixture model is not identifiable because the number of moments is one and is not larger than the number of parameters. In a univariate case, if $F_1 = p^{(1)}$, the PMF becomes

$$f(y_1) = p^{(1)}I(y_1 = 0) + (1 - p^{(1)}) [1 - p^{(1)}]^{1-y_1} [p^{(1)}]^{y_1}. \quad (12)$$

Then, the parameter $p^{(1)}$ is identifiable but this model is not practical in applied econometric fields.

3. Monte Carlo simulation results

This section presents results of Monte Carlo simulation of a bivariate zero-inflated binary model. The sample sizes used are 500 and 2,000, and the number of simulations in all experiments is set at 1,000. We generate two binary variables y_1 and y_2 with probabilities of F_1 and F_2 , respectively. To simplify the calculation, we consider one parameter model and do not include any regressors (covariates). The proportion of these two variables is $1 - p^{(1)}$ and the rest $p^{(1)}$ are zeros.

Table I presents true values of $(F_1, F_2, p^{(1)})$ and the Monte Carlo simulation results of the bivariate zero-inflated binomial model. The results of $(F_1, F_2, p^{(1)})$ show that the parameter estimates are unbiased. The root mean squared errors (RMSE) decrease when the sample size increases in each simulation. Moreover, if $p^{(1)}$ is large ($1 - p^{(1)}$ is small), the bias for (F_1, F_2) is large since the non-zero value is small.

4. An application in the RAND Health Insurance Experiment

Using the same data from the RAND Health Insurance Experiment (RAND HIE) of Deb and Trivedi (2002) and Cameron and Trivedi (2010, Ch.15), we

Table I: Monte Carlo simulation results

	True	$N = 500$	$N = 2,000$	True	$N = 500$	$N = 2,000$
F_1	0.25	0.001 (0.034)	0.001 (0.017)	0.25	0.005 (0.089)	0.001 (0.042)
F_2	0.5	0.000 (0.055)	0.001 (0.027)	0.5	0.005 (0.146)	0.004 (0.067)
$p^{(1)}$	0.3	-0.006 (0.072)	0.000 (0.034)	0.9	-0.008 (0.048)	-0.001 (0.014)

	True	$N = 500$	$N = 2,000$	True	$N = 500$	$N = 2,000$
F_1	0.2	0.002 (0.039)	0.001 (0.020)	0.2	0.008 (0.108)	0.002 (0.053)
F_2	0.3	0.003 (0.054)	0.001 (0.028)	0.3	0.011 (0.152)	0.002 (0.073)
$p^{(1)}$	0.3	-0.013 (0.117)	-0.002 (0.060)	0.9	-0.053 (0.193)	-0.006 (0.028)

Note: The mean bias reports appear without parentheses. RMSE are in parentheses.

analyze the difference between a bivariate zero-inflated binomial model and a bivariate binary probit model without inflation. We use `year = 3` and female samples from this data. The number of observations is $N = 2,875$. The first binary outcome is `DMENTVIS` (its mean is 4.28%) for an individual visiting psychotherapy, and the second is `DNOTMD` (its mean is 20.52%) for an individual visiting a non-medical doctor in the current year. When analyzing these outcomes, a bivariate binary probit model is usually applied. However, since there exist many zero-valued observations, a zero-inflated model, such as a bivariate zero-inflated binomial model discussed in Section 2, may be more suitable to analyze these bivariate binary outcomes. The regressors are age (`AGE`; its average is 26.97), log of family income (`LINC`; its average is 8.65), and the number of chronic diseases (`NDISEASE`; its average is 12.52).

In a bivariate zero-inflated binomial model, we specify both F_1 and F_2 as probit models, $\Phi(\mathbf{x}'\boldsymbol{\beta}_1)$ and $\Phi(\mathbf{x}'\boldsymbol{\beta}_2)$. The $\Phi(\cdot)$ is a cumulative distribution of a standard normal distribution, $\boldsymbol{\beta}_1$ and $\boldsymbol{\beta}_2$ are parameter vectors to be estimated, and \mathbf{x} is a vector of regressors. Moreover, a component proportion $p^{(1)}$ is $1 - \Phi(\mathbf{x}'\boldsymbol{\beta}_p)$. An increase in \mathbf{x} means a more likely transition from 0 to 0/1. Then, the PMF of this model is

$$f(y_1, y_2) = (1 - \Phi(\mathbf{x}'\boldsymbol{\beta}_p)) I(y_1 = y_2 = 0) \\ + \Phi(\mathbf{x}'\boldsymbol{\beta}_p) [1 - \Phi(\mathbf{x}'\boldsymbol{\beta}_1)]^{1-y_1} [\Phi(\mathbf{x}'\boldsymbol{\beta}_1)]^{y_1} [1 - \Phi(\mathbf{x}'\boldsymbol{\beta}_2)]^{1-y_2} [\Phi(\mathbf{x}'\boldsymbol{\beta}_2)]^{y_2}.$$

This PMF is essentially the same as (11) and is identifiable from Proposition 2. Table II displays the estimated results of the two models, a bivariate zero-inflated binomial model and a bivariate binary probit model, as well as values of

Table II: Estimated results of RAND HIE

	Zero-inflated			Bivariate probit		
DMENTVIS						
AGE	0.011	(0.003)	***	0.008	(0.003)	***
LINC	0.017	(0.104)		0.056	(0.042)	
NDISEASE	0.014	(0.006)	**	0.013	(0.006)	**
constant	-2.331	(0.892)	***	-2.603	(0.382)	***
DNOTMD						
AGE	0.019	(0.002)	***	0.013	(0.002)	***
LINC	-0.050	(0.028)	*	0.024	(0.023)	
NDISEASE	0.016	(0.004)	***	0.014	(0.004)	***
constant	-1.081	(0.246)	***	-1.594	(0.204)	***
$\Phi(\cdot) = 1 - p^{(1)}$						
AGE	-0.066	(0.020)	***			
LINC	0.362	(0.099)	***			
NDISEASE	-0.002	(0.015)				
constant	1.260	(0.719)	*			
ρ				0.126	(0.056)	**
log-likelihood	-1,887.012			-1,896.973		
AIC	3,798.024			3,811.945		
BIC	3,869.589			3,865.620		

Notes: Standard errors are in parentheses; Statistically significant at the 1% (***), 5% (**), and 10% (*) levels; $AIC = -2 \ln L + 2K_L$, $BIC = -2 \ln L + K_L \ln N$, where L is the maximum likelihood, K_L is the number of parameters, and N is the number of observations ($N = 2,875$).

the log-likelihood, Akaike's information criteria (AIC), and Bayesian information criteria (BIC). The maximum value of the log-likelihood and the minimum value of the AIC determine the bivariate zero-inflated binomial model. The minimum value of the BIC determines the bivariate binary probit model.

In Table II, we observe some features of the estimated parameters. First, the estimated parameters of the two models resemble each other except for the constant terms and insignificant variables, such as LINC. Even though the bivariate binary probit model allows the two variables to be correlated and the zero-inflated model assumes independence of two binary variables, the zeros capture the correlation. Second, the significance levels of the estimated parameters resemble each other. Third, in the transition from 0 to 0/1 of the zero-inflated model, the log of income is statistically significant at the 1% level.

This indicates a difference in the interpretation of marginal effects of regressors.

Based on estimated results in Table II, we calculate marginal effects of regressors. The average marginal effect of AGE to DMENTVIS (visiting psychotherapy) in the bivariate probit model is 0.04% points, statistically significant at the 1% level. In the bivariate zero-inflated binomial model, those values are -0.04% points from a non-user to a user and 0.1% points on 0/1 decision making. Both values are statistically significant at the 1% level. In the bivariate probit model, however, the LINC has no effect on DMENTVIS (the average marginal effect is 0.2% points but not statistically significant), its value is 0.2% points from a non-user to a user and statistically significant in the bivariate zero-inflated binomial model (but has no effect on the 0/1 decision making).

5. Conclusions

This study analyzes identifiability of finite mixture models with the existence of an MGF. Our results show that a finite mixture model is identifiable when the Jacobian determinant of the joint moments and sample moments is not zero. This paper obtains the conditions of a zero-inflated binomial model, which has the simplified structure of a finite mixture model, and demonstrates that a bivariate zero-inflated binomial model is identifiable. Monte Carlo experiments support our demonstration and show good performance. Using the RAND HIE data, we compare the differences of estimated coefficients between a bivariate zero-inflated binomial model and a bivariate binary probit model without inflation. We decompose the marginal effects of moving a non-user to a user and the 0/1 decision making of users. Based on the bivariate zero-inflated binomial model, increasing family income slightly converts a non-user into a user, but has no effect in the bivariate binary probit model. These results suggest the usefulness of the finite mixture models. When analyzing multi-variate binary variables including panel data, it is feasible to estimate finite mixture models using Proposition 1. If the estimated result of some variables is not significant, finite mixture models are good alternatives.

References

- Blischke, W.R. (1964) "Estimating the Parameters of Mixtures of Binomial Distributions" *Journal of the American Statistical Association* **59 (306)**, 510–528.
- Cameron, A.C. and P.K. Trivedi (2010) *Microeconometrics Using Stata, Revised Edition*, Stata Press: Texas.
- Deb, P. and P.K. Trivedi (1997) "Demand for Medical Care by the Elderly: A Finite Mixture Approach" *Journal of Applied Econometrics* **12 (3)**, 313–336.

- Deb, P. and P.K. Trivedi (2002) “The Structure of Demand for Health Care: Latent Class Versus Two-Part Models” *Journal of Health Economics* **21** (4), 601–625.
- Kasahara, H. and K. Shimotsu (2014) “Non-parametric Identification and Estimation of the Number of Components in Multivariate Mixtures” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76** (1), 97–111.
- Teicher, H. (1963) “Identifiability of Finite Mixtures” *Annals of Mathematical Statistics* **34** (4), 1265–1269.
- van Wieringen, W.N. (2005) “On Identifiability of Certain Latent Class Models” *Statistics & Probability Letters* **75** (3), 211–218.
- Winkelmann, R. (2004) “Health Care Reform and the Number of Doctor Visits; An Econometric Analysis” *Journal of Applied Econometrics* **19** (4), 455–472.