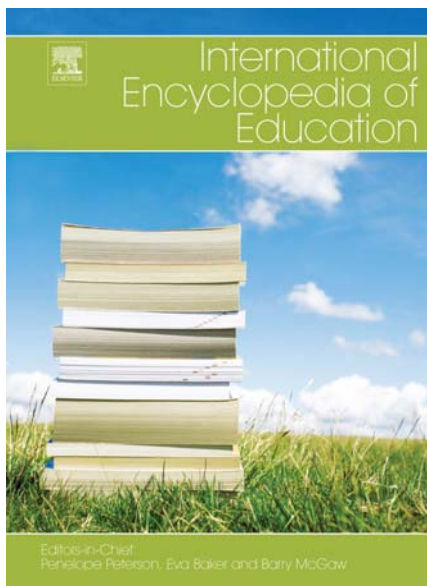


**Provided for non-commercial research and educational use.  
Not for reproduction, distribution or commercial use.**

This article was originally published in the *International Encyclopedia of Education* published by Elsevier, and the attached copy is provided by Elsevier for the author's benefit and for the benefit of the author's institution, for non-commercial research and educational use including without limitation use in instruction at your institution, sending it to specific colleagues who you know, and providing a copy to your institution's administrator.



All other uses, reproduction and distribution, including without limitation commercial reprints, selling or licensing copies or access, or posting on open internet sites, your personal or institution's website or repository, are prohibited. For exceptions, permission may be sought for such use through Elsevier's permissions site at:

<http://www.elsevier.com/locate/permissionusematerial>

Vermunt J K (2010), Latent Class Models. In: Penelope Peterson, Eva Baker, Barry McGaw, (Editors), *International Encyclopedia of Education*. volume 7, pp. 238-244. Oxford: Elsevier.

## Latent Class Models

J K Vermunt, Tilburg University, Tilburg, The Netherlands

© 2010 Elsevier Ltd. All rights reserved.

### Introduction

A statistical model can be called a latent class (LC) or mixture model if it assumes that some of its parameters differ across unobserved subgroups, LCs, or mixture components. This rather general idea has several seemingly unrelated applications, the most important of which are clustering, scaling, density estimation, and random-effects modeling. It should be noted that in applied fields, the terms LC model and mixture model are often used interchangeably, which is also what is done here. In the more technical statistical literature on mixture modeling, the term LC analysis is reserved for a specific type of mixture model, that is, a mixture model for a set of categorical items (for the classical LC model).

LC analysis was introduced in 1950 by Lazarsfeld as a tool for building typologies (or clustering) based on dichotomous observed variables (Lazarsfeld, 1950). More than 20 years later, Goodman (1974) made this model applicable in practice by developing an algorithm for obtaining maximum-likelihood estimates of the model parameters, as well as proposed extensions for polytomous manifest variables and did important work on the issue of model identification. Many important extensions of this classical LC model have been proposed since then, such as models containing explanatory variables (Dayton and Macready, 1988), models that relax the local-independence assumption (Hagenaars, 1988), constrained models similar to item response theory (IRT) models (Lindsay *et al.*, 1991; Heinen, 1996), models with multiple latent variables (Magidson and Vermunt, 2001), models for longitudinal data (Van de Pol and Langeheine, 1990), and models for multilevel data (Vermunt, 2003).

Whereas this classical LC model and its extensions are conceived primarily as a clustering and scaling tool for categorical data analysis, LC and finite-mixture models can be useful in several other areas as well. One of these is as a probabilistic cluster-analysis tool for continuous observed variables, an approach that offers many advantages over traditional cluster techniques such as K-means clustering (Wolfe, 1970; McLachlan and Peel, 2000; Vermunt and Magidson, 2002). Another application area is dealing with unobserved heterogeneity, as happens in mixture regression analysis of multilevel or repeated-measurement data (Wedel and DeSarbo, 1994; Vermunt and Van Dijk, 2001).

The remainder of this article is organized as follows. After introducing the simplest type of LC models, various

restricted LC models as well as models with explanatory variables are discussed. Next, an overview of other types of LC and mixture models, which includes various recently proposed extensions, is presented. In the end, attention is paid to parameter estimation, model selection, and software.

### Simple LC and Mixture Models

LC analysis is typically used as a tool for analyzing multivariate response data; that is, data consisting of several dependent variables, response variables, or items. We denote the response of subject  $i$  on dependent variable  $j$  by  $y_{ij}$ , and the number of dependent variables by  $\mathcal{J}$ . The full-response vector of a subject is denoted by  $\mathbf{y}_i$ . To make things more concrete, **Table 1** presents a small illustrative data set consisting of three dichotomous responses,  $y_{i1}$ ,  $y_{i2}$ , and  $y_{i3}$  (0 = incorrect; 1 = correct). This is a subset of items from a mathematics test administered to 2156 children. The frequency column contains the observed frequency count for each of the eight possible answer patterns.

In addition to the  $\mathcal{J}$  observed dependent variables, an LC model contains a discrete latent variable. We denote a subject's unobserved score on this latent variable by  $v_i$ , the number of LCs by  $C$ , and a particular class by  $c$ , where  $c = 1, 2, \dots, C$ . The aim of an LC analysis of the data set in **Table 1** could be to classify pupils into two groups, masters and nonmasters, which differ with respect to the probability of answering the test items correctly. The results obtained with a two-class model will be used to illustrate the various components of an LC model.

LC analysis defines a model for  $f(\mathbf{y}_i)$ , the probability density of the multivariate response vector  $\mathbf{y}_i$ . In the above example, this is the probability of answering the items according to one of the eight possible response patterns, for example, of answering the first two items correctly and the last one incorrectly, which as can be seen in **Table 1** equals 0.161 for the estimated two-class model. The assumption underlying any type of LC or mixture model is that the density  $f(\mathbf{y}_i)$  is a weighted average (or mixture) of the  $C$  class-specific densities  $f(\mathbf{y}_i|v_i = c)$ . This is expressed mathematically as follows:

$$f(\mathbf{y}_i) = \sum_{c=1}^C P(v_i = c) f(\mathbf{y}_i|v_i = c). \quad [1]$$

Here,  $P(v_i = c)$  denotes the probability that a subject belongs to LC  $c$ . For our small empirical example, the estimates of these (prior) class membership probabilities

**Table 1** Small data set with three dichotomous responses

$y_{i1}$	$y_{i2}$	$y_{i3}$	Frequency	$f(y_i v_i = 1)$	$f(y_i v_i = 2)$	$f(y_i)$	$f(v_i = 1 y_i)$	$f(v_i = 2 y_i)$	Modal
0	0	0	239	0.004	0.272	0.111	0.020	0.980	2
0	0	1	101	0.010	0.102	0.047	0.128	0.872	2
0	1	0	283	0.038	0.271	0.131	0.175	0.825	2
0	1	1	222	0.104	0.102	0.103	0.605	0.395	1
1	0	0	105	0.020	0.092	0.049	0.248	0.753	2
1	0	1	100	0.054	0.035	0.046	0.703	0.297	1
1	1	0	348	0.208	0.091	0.161	0.774	0.226	1
1	1	1	758	0.562	0.034	0.352	0.961	0.039	1

are 0.601 and 0.399 for classes 1 and 2, respectively (see **Table 2**). The assumed mechanism by eqn [1] is that each individual belongs to one of  $C$  exhaustive and mutually exclusive classes with probability  $P(v_i = c)$  and that given membership of LC  $c$  one provides responses according to the probability density associated to this class. **Table 1** shows the estimated values of  $f(y_i|v_i = 1)$  and  $f(y_i|v_i = 2)$  for our data sets. As can be seen, LC 1 has higher probabilities for the response patterns with 2 or 3 correctly answered items, whereas class 2 has higher probabilities for the response patterns with 0 or 1 item correct.

The classical LC model combines the assumption of eqn [1] shared by all mixture models with the assumption of local independence. Local independence means that the  $\mathcal{J}$  responses are mutually independent given a subject's class membership. It can be expressed as follows:

$$f(y_i|v_i = c) = \prod_{j=1}^{\mathcal{J}} f(y_{ij}|v_i = c). \quad [2]$$

Independence implies that the joint density  $f(y_i|v_i = c)$  is obtained as a product of the  $\mathcal{J}$  item-specific densities  $f(y_{ij}|v_i = c)$ . In our example  $f(y_{ij} = 1|v_i = c)$  is the class-specific probability of giving a correct answer to item  $j$ . As reported in **Table 2**, for a subject belonging to the first LC, these equal 0.844, 0.912, and 0.730 for items 1, 2, and 3, respectively. The local independence assumption implies, for example, that the probability of answering the first two items correctly and the last one incorrectly for someone in LC one equals  $0.844 \times 0.912 \times (1-0.730) = 0.208$ . Note that the local independence assumption is also used in other types of latent variables models, such as in factor analysis and IRT modeling, and is thus not specific for LC analysis.

Combining the two basic eqns [1] and [2] yields the following model for  $f(y_i)$ :

$$f(y_i) = \sum_{c=1}^C P(v_i = c) \prod_{j=1}^{\mathcal{J}} f(y_{ij}|v_i = c). \quad [3]$$

To complete the model specification, we need to define the form of the conditional densities  $f(y_{ij}|v_i = c)$ . In the classical LC model for categorical items, these are multinomial probability densities; that is,

**Table 2** Parameters (class proportions and probability of a correct answer) obtained with two-class model for data in **Table 1**

	$c = 1$	$c = 2$
$P(v_i = c)$	0.601	0.399
$\pi_{11c} = P(y_{i1} = 1 v_i = c)$	0.844	0.252
$\pi_{21c} = P(y_{i2} = 1 v_i = c)$	0.912	0.499
$\pi_{31c} = P(y_{i3} = 1 v_i = c)$	0.730	0.273

$$f(y_{ij}|v_i = c) = \prod_{r=0}^{R_j-1} \pi_{jrc}^{y_{ijr}^*},$$

where  $R_j$  is the number of categories of item  $j$ ,  $0 \leq y_{ij} \leq R_j - 1$ , and  $y_{ijr}^* = 1$  if  $y_{ij} = r$  and 0 otherwise. Note this is a slightly complicated, but mathematically elegant, way to express that someone in LC  $c$  has a probability equal to  $\pi_{jrc} = P(y_{ij} = r|v_i = c)$  of giving response  $r$  to item  $j$ . In the special case of a dichotomous response, the multinomial distribution reduces to the Bernoulli distribution with success probability  $\pi_{jc} = \pi_{j1c} = P(y_{ij} = 1|v_i = c)$ . **Table 2** presents these probabilities for our small example.

It is important to note that LC models cannot only be used with categorical responses, but also with continuous responses and counts. The density  $f(y_{ij}|v_i = c)$  could be a binomial, Poisson, or negative binomial distribution for counts, and a normal or gamma distribution for continuous responses. The mixture model for continuous response variables is sometimes referred to as the latent profile model. The parameters of this model are the class proportions and class-specific item means and variances ( $\mu_{jc}$  and  $\sigma_{jc}^2$ ).

By comparing the  $\mathcal{J}$  sets of item parameters across classes, one can name the classes. The parameter estimates presented in **Table 2** show that the first class can be named the masters because pupils belonging to that class have much higher probabilities of answering the test items correctly than pupils belonging to the second non-masters class.

Similar to cluster analysis, one of the purposes of LC analysis might be to assign individuals to LCs.

The probability of belonging to LC  $c$  given responses  $\mathbf{y}_i$  – often referred to as posterior membership probability – can be obtained by the Bayes rule:

$$P(v_i = c | \mathbf{y}_i) = \frac{P(v_i = c) f(\mathbf{y}_i | v_i = c)}{f(\mathbf{y}_i)} \quad [4]$$

**Table 1** reports  $P(v_i = c | \mathbf{y}_i)$  for each answer pattern. For example,  $P(v_i = 1 | \mathbf{y}_i)$  equals 0.774 for the (1,1,0) pattern, which is obtained as  $0.601 \times 0.208 / 0.161$ . The most common classification rule is modal assignment, which amounts to assigning each individual to the LC with the highest  $f(v_i = c | \mathbf{y}_i)$ . The last column of **Table 1** reporting the modal assignments shows that pupils with at least two correct answers are assigned to class 1 and the others to class 2.

In the introduction, we stated that mixture models are statistical models in which parameters are assumed to differ across LCs. But what is the statistical model used in the simple LC models discussed so far? It is the independence model: we assume responses to be independent, with different parameter values for each class. Depending on the scale type of the response variables, these parameters are Bernoulli probabilities, multinomial probabilities, normal means and variances, Poisson rates, etc.

### Generalized Linear Models for Item Probabilities/Mean

**Haberman (1979)** showed that the LC model for categorical response variables can also be specified as a log-linear model for an expanded table, including the latent variable  $v_i$  as an additional dimension. Using such a log-linear specification is equivalent to parameterizing the response probability for item  $j$  as follows:

$$\log\left(\frac{\pi_{jrc}}{\pi_{j0c}}\right) = \log\left(\frac{P(\mathcal{Y}_{ij} = r | v_i = c)}{P(\mathcal{Y}_{ij} = 0 | v_i = c)}\right) = \alpha_{jr} + \beta_{jcr}, \quad [5]$$

for  $1 \leq r \leq R_j - 1$ ; that is, as a multinomial logistic regression models with intercepts  $\alpha_{jr}$  and slopes  $\beta_{jcr}$  (note that we use the first item category,  $r = 0$ , as baseline). One identification constraint needs to be imposed, for example,  $\beta_{j1r} = 0$  (the parameters for class 1 are fixed to 0) or  $\alpha_{jr} = 0$  (intercepts are fixed to 0).

For dichotomous responses and binomial counts, the regression model could be a binary logit or probit model, for Poisson counts a log-linear model, and for continuous responses a standard linear model. These are generalized linear models (GLMs) of the form

$$g[E(\mathcal{Y}_{ij} | v_i = c)] = \alpha_j + \beta_{jc}, \quad [6]$$

where  $g[\cdot]$  is the link function transforming the expected value of  $\mathcal{Y}_{ij}$  to the linear term. For ordinal polytomous variables, one may use an ordinal regression model,

such as an adjacent-category or cumulative logit model. These are models that restrict the item response probabilities  $\pi_{jrc}$ .

### Some Restricted Models for Categorical Items

Many interesting types of restricted LC models for categorical items have been proposed, which involve imposing (linear) constraints on either the conditional probabilities  $\pi_{jrc}$  or the logit coefficients of eqn [5]. One of these is the probabilistic Guttman scaling model for dichotomous responses, which is an LC model with  $C = \mathcal{F} + 1$  classes, one for each possible total score. The idea is that apart from measurement error, class  $c$  should provide a negative answer to the  $c - 1$  easiest items and a positive answer to the remaining  $\mathcal{F} - (c - 1)$  items. The various types of probabilistic Guttman models differ in the constraints they impose on the measurement error. The simplest and most restricted model is the **Proctor (1970)** model.

**Table 3** presents the parameter estimates obtained when fitting the Proctor model to the data set in **Table 1**. As can be seen, the probability of a correct response is either 0.833 or  $0.167 = 1 - 0.833$ . The measurement error – or the probability of giving a response which is not in agreement with the class – is estimated to be equal to 0.167. Whereas the Proctor model assumes that the measurement error is constant across items and classes, less-restricted models can be defined which allow the error probabilities to differ across items, classes, or both (see, e.g., **Dayton, 1999**). Note that these equality constraints on the error probabilities can also be defined using linear constraints on the logit parameters:  $\alpha_{j1} = 0$  and  $\beta_{j1c} = -\beta^*$  for  $c \leq j$  and  $\beta_{jc} = \beta^*$  otherwise.

**Croon (1990)** proposed a restricted LC model that similar to nonparametric IRT (**Sijtsma and Molenaar, 2000**) assumes monotonic item response functions; that is,  $\pi_{j1c} \leq \pi_{j1,c+1}$ , or, equivalently,  $\beta_{j1c} \leq \beta_{j1,c+1}$ . A more restricted version, in which not only classes but also items are ordered, is obtained by imposing the additional set of restriction  $\pi_{j+11c} \leq \pi_{j1c}$ ; that is, by assuming double monotony. **Vermunt (2001)** has discussed various generalizations of these models.

**Table 3** Parameters (class proportions and probability of a correct answer) obtained with proctor model for data in **Table 1**

	$c = 1$	$c = 2$	$c = 3$	$c = 4$
$P(v_i = c)$	0.160	0.155	0.126	0.559
$\pi_{21c} = P(\mathcal{Y}_{i2} = 1   v_i = c)$	0.167	0.833	0.833	0.833
$\pi_{11c} = P(\mathcal{Y}_{i1} = 1   v_i = c)$	0.167	0.167	0.833	0.833
$\pi_{31c} = P(\mathcal{Y}_{i3} = 1   v_i = c)$	0.167	0.167	0.167	0.833

Various authors described the connection between restricted LC analysis and parametric IRT modeling (see, e.g., Heinen, 1996; Lindsay *et al.*, 1991); that is, IRT models with a discrete specification of the distribution of the underlying trait or ability can be defined as LC models with restrictions on the logistic parameters. The key restriction is  $\beta_{jrc} = \beta_{jr}^* \cdot \theta_c$  for nominal items and  $\beta_{jrc} = \beta_j^* \cdot r \cdot \theta_c$  for ordinal items, where  $\theta_c$  are LC locations representing the  $C$  possible values of the discretized latent trait. These locations may be fixed *a priori*, for example, at  $-2, -1, 0, 1,$  and  $2$  in the case of  $C = 5$ , but may also be treated as free parameters to be estimated. Depending on whether the items are dichotomous, ordinal, or nominal, this yields a 2-parameter logistic, generalized partial credit, or nominal response model. Further restrictions involve equating  $\beta_j^*$  across items, yielding Rasch and partial credit models, and imposing across-category and across-item restrictions on  $\alpha_{jr}$  parameters as in rating scale models for ordinal items.

## Models with Explanatory Variables

The most important extension of the LC models discussed so far is the possibility to include explanatory variables affecting the responses (Wedel and DeSarbo, 1994) or the class memberships (Dayton and Macready, 1988). Denoting the vector with explanatory variables for subject  $i$  by  $\mathbf{x}_i$ , the LC model of interest can be formulated as follows:

$$f(\mathbf{y}_i | \mathbf{x}_i) = \sum_{c=1}^C P(v_i = c | \mathbf{x}_i) \prod_{j=1}^{j_i} f(y_{ij} | v_i = c, \mathbf{x}_i). \quad [7]$$

The main difference compared to the model defined in eqn [3] is that now we have a model for  $f(\mathbf{y}_i | \mathbf{x}_i)$  – the conditional density of  $\mathbf{y}_i$  given  $\mathbf{x}_i$ .

Similar to the regression models for the response variables introduced in eqns [5] and [6], one can define a mixture regression model with explanatory variables; that is,

$$g[E(y_{ij} | v_i = c, \mathbf{x}_{ij})] = \alpha_c + \sum_{p=1}^P \beta_{pc} x_{ijp}. \quad [8]$$

As before,  $y_{ij}$  may refer to the response on item  $j$  by pupil  $i$ , in which case the explanatory variables will consist of a design matrix defining the item parameters. However, the model in eqn [8] can also be used for many other purposes. In fact, it is a model for analyzing two-level data sets, where regression parameters are allowed to differ across LCs (of higher-level units). For example,  $y_{ij}$  could be the test score of pupil  $j$  belonging to school  $i$ , and  $\mathbf{x}_{ij}$  a set of pupil characteristics (e.g., intelligence quotient (IQ)). A mixture regression model would identify LCs of schools with different intercepts and different effects of child characteristics on the test scores.

Another possible application is in the analysis of longitudinal data, where  $j$  is a time point for subject  $i$ , and where vector  $\mathbf{x}_{ij}$  contains time variables. This yields a LC growth model in which subjects are grouped based on their developmental trajectories (Vermunt, 2007). A last possible application that can be mentioned is in experiments in which subjects are observed in multiple conditions, such as in conjoint studies. The mixture regression model can be used to group subjects based on their reactions on the experimental conditions. In fact, in each of these application types, the LC model is used as a random-coefficient model without parametric assumptions about the distribution of the random effects (Aitkin, 1999; Vermunt and Van Dijk, 2001).

As shown in eqn [7], an individual's class membership may also be predicted using covariates. This is achieved by defining a multinomial logistic regression model for  $P(v_i = c | \mathbf{x}_i)$ :

$$\log \frac{P(v_i = c | \mathbf{x}_i)}{P(v_i = 1 | \mathbf{x}_i)} = \gamma_{0c} + \sum_{q=1}^Q \gamma_{qc} x_{iq}.$$

Strongly related are multiple-group LC models. These can be defined using the grouping variable as a nominal explanatory in the model.

## Extensions

The most common model-fitting strategy in LC analysis is to increase the number of classes until the local-independence assumption holds. This may, however, yield solutions which are difficult to interpret. One alternative approach is to relax the local-independence assumption by allowing for associations between particular item pairs. Hagnaars (1988) showed how to define LC models with local dependencies for categorical responses. With continuous responses, this is easily achieved using multivariate instead of univariate normal distributions for locally dependent items (see, e.g., McLachlan and Peel, 2000; Vermunt and Magidson, 2002).

Another alternative strategy involves increasing the number of discrete latent variables instead of the number of LCs, which is especially useful if the items measure several dimensions. This so-called discrete-factor modeling approach (Magidson and Vermunt, 2001) is a special case of the path-modeling approach for discrete latent variables developed by Hagnaars (1990) and Vermunt (1997). Many other interesting models can be defined within this framework, such as latent Markov models for the analysis of longitudinal data (Van de Pol and Langeheine, 1990) and LC models for cognitive diagnosis (De la Torre and Douglas, 2004).

Recently, models have been developed that contain both discrete and continuous latent variables. Examples of these are mixture factor models (Yung, 1997; McLachlan



and Peel, 2000), mixture structural equation models (Dolan and Van der Maas, 1997), and mixture IRT models (Rost, 1990).

Probably the most recent extension is the multilevel LC model (Vermunt, 2003). One of its variants is a model with discrete latent variables at multiple levels of a hierarchical structure: for example, children belong to LCs with different performances on a set of test items, and schools belong to LCs with different distributions of children across the child-level performance classes. Multilevel LC models can be used for the analysis of two-level multivariate and three-level univariate response data.

## Maximum Likelihood Estimation

The parameters of LC models are typically estimated by means of maximum likelihood (ML). The log-likelihood function that is maximized is based on the probability densities defined in eqns [1–3]; that is,

$$\ln L = \sum_{i=1}^N \ln f(\mathbf{y}_i).$$

With categorical responses one will typically group the data and construct a frequency table as we did in Table 1. The log-likelihood function for grouped data equals

$$\ln L = \sum_{k=1}^K n_k \ln f(\mathbf{y}_k),$$

where  $k$  is a data pattern,  $K$  the number of different data patterns, and  $n_k$  the cell count corresponding to data pattern  $k$ . Notice that only nonzero observed cell entries contribute to the log-likelihood function, a feature that is exploited by several more efficient LC software packages that have been developed within the past few years.

One of the problems in the estimation of LC models for discrete  $y_{ij}$  is that model parameters may be nonidentified, even if the number of degrees of freedom – the number of independent cells in the  $\mathcal{F}$ -way cross-tabulation minus the number of free parameters – is larger or equal to zero. Nonidentification means that different sets of parameter values yield the same maximum of the log-likelihood function or, worded differently, that there is no unique set of parameter estimates. The formal identification check is via the Jacobian matrix (matrix of first derivatives of  $f(\mathbf{y}_i)$ ), which should be column full rank. Another option is to estimate the model of interest with different sets of starting values. Except for local solutions (see below), an identified model gives the same final estimates for each set of the starting values.

Although there are no general rules with respect to the identification of LC models, it is possible to provide certain minimal requirements and point to possible pitfalls. For an unrestricted LC analysis, one needs at least

three responses ( $y_{ij}$ 's) per individual, but if these are dichotomous, no more than two LCs can be identified. One has to watch out with four dichotomous response variables, in which case the unrestricted three-class model is not identified, even though it has a positive number of degrees of freedom. With five dichotomous items, however, even a five-class model is identified. Usually, it is possible to achieve identification by constraining certain model parameters.

A second problem associated with the estimation of LC models is the presence of local maxima. The log-likelihood function of an LC model is not always concave, which means that hill-climbing algorithms may converge to a different maximum depending on the starting values. Usually, we are looking for the global maximum. The best way to proceed is, therefore, to estimate the model with different sets of random starting values. Typically, several sets converge to the same highest log-likelihood value, which can then be assumed to be the ML solution. Some software packages have automated the use of multiple sets of random starting values to reduce the probability of getting a local solution.

Another problem in LC modeling is the occurrence of boundary solutions, which are probabilities equal to 0 (or 1) or logit parameters equal to minus (or plus) infinity. These may cause numerical problems in the estimation algorithms, occurrence of local solutions, and complications in the computation of standard errors and number of degrees of freedom of the goodness-of-fit tests. Boundary solutions can be prevented by imposing constraints or by taking into account other kinds of prior information on the model parameters.

The most popular methods for solving the ML estimation problem are the expectation–maximization (EM) and Newton–Raphson (NR) algorithms. EM is a very stable iterative method for ML estimation with incomplete data. NR is a faster procedure that, however, needs good starting values to converge. The latter method makes use of the matrix of second-order derivatives of the log-likelihood function, which is also needed for obtaining standard errors of the model parameters.

## Model Selection Issues

The goodness-of-fit of LC models for categorical responses can be tested using Pearson and likelihood-ratio chi-squared tests. The latter is defined as

$$L^2 = 2 \sum_{k=1}^K n_k \ln \frac{n_k}{N \cdot f(\mathbf{y}_k)}.$$

As in log-linear analysis, the number of degrees of freedom (df) equals the number of cells in the frequency table minus 1, minus the number of independent parameters. In an unrestricted LC model,

$$df = \prod_{j=1}^{\mathcal{F}} R_j - C \cdot \left[ 1 + \sum_{j=1}^{\mathcal{F}} (R_j - 1) \right].$$

Although it is no problem to estimate LC models with 10, 20, or 50 indicators, in such cases, the frequency table may become very sparse and, as a result, asymptotic  $p$ -values can no longer be trusted. An elegant but somewhat time-consuming solution to this problem is to estimate the  $p$ -values by parametric bootstrapping. Another option is to assess model fit in lower order marginal tables (e.g., in the two-way marginal tables).

Even though models with  $C$  and  $C + 1$  are nested, one cannot test them against each other using a standard likelihood-ratio test because it does not have an asymptotic chi-squared distribution. A way out to this problem is to approximate its sampling distribution using bootstrapping. But since this method is computationally demanding, usually alternative methods are required for comparing models with different numbers of classes. One popular method is the use of information criteria such as the Bayesian information criterion (BIC) and Akaike information criterion (AIC). Another more descriptive method is a measure for the proportion of total association accounted for by a  $C$ -class model,  $[L^2(1) - L^2(C)]/L^2(1)$ , where the  $L^2$  value of the one-class (independence) model,  $L^2(1)$ , is used as a measure of total association in the  $\mathcal{F}$ -way frequency table.

Usually, we are not only interested in goodness-of-fit but also in the performance of the modal classification rule (see eqn [4]). The estimated proportion of classification errors under modal classification equals

$$E = \sum_{i=1}^N \frac{1}{N} \{1 - \max[P(v_i = c | \mathbf{y}_i)]\}.$$

This number can be compared to the proportion of classification errors based on the unconditional probabilities  $P(v_i = c)$ , yielding a reduction of errors measure

$$\lambda = 1 - \frac{E}{1 - \max[P(v_i = c)]}.$$

The closer this nominal  $R^2$ -type measure is to 1, the better the classification performance of a model. Other types of classification error-reduction measures have been proposed based on entropy or qualitative variance.

## Software

One of the first LC analysis programs, maximum likelihood latent structure analysis (MLLSA), made available by Clifford Clogg in 1977, was limited to a relatively small number of nominal variables. Today's programs can handle many more variables, as well as other scale types. For example, the LEM program (Vermunt, 1997) provides

a command language that can be used to specify a large variety of models for categorical data, including LC models. Mplus is a command-language-based structural-equation modeling package that implements many types of LC and mixture models. In addition, routines for the estimation of specific types of LC models are available as SAS, R, and Stata macros (see, e.g., Lanza *et al.*, 2007; Skrondal and Rabe-Hesketh, 2004).

Latent GOLD is a program that was especially developed for LC analysis, and which contains both an SPSS-like point and click-user interface and a syntax language. It implements all important types of LC models, such as models for response variables of different scale types, restricted LC models, models with predictors, models with local dependencies, models with multiple discrete latent variables, LC path models, LC Markov models, mixture factor analysis and IRT, and multilevel LC models, as well as features for dealing with partially missing data, for performing bootstrapping, and for dealing with complex samples.

## Bibliography

- Aitkin, M. (1999). A general maximum likelihood analysis of variance components in generalized linear models. *Biometrics* **55**, 218–234.
- Croon, M. A. (1990). Latent class analysis with ordered latent classes. *British Journal of Mathematical and Statistical Psychology* **43**, 171–192.
- Dayton, C. M. (1999). *Latent Class Scaling Analysis*. Thousand Oaks, CA: Sage.
- Dayton, C. M. and Macready, G. B. (1988). Concomitant-variable latent class models. *Journal of the American Statistical Association* **83**, 173–178.
- De la Torre, J. and Douglas, J. (2004). Higher-order latent trait models for cognitive diagnosis. *Psychometrika* **69**, 333–353.
- Dolan, C. V. and Van der Maas, H. L. J. (1997). Fitting multivariate normal finite mixtures subject to structural equation modeling. *Psychometrika* **63**, 227–253.
- Goodman, L. A. (1974). Exploratory latent structure analysis using both identifiable and unidentifiable models. *Biometrika* **61**, 215–231.
- Haberman, S. J. (1979). *Analysis of Qualitative Data. Vol. 2: New Developments*. New York: Academic Press.
- Hagenaars, J. A. (1988). Latent structure models with direct effects between indicators: Local dependence models. *Sociological Methods and Research* **16**, 379–405.
- Hagenaars, J. A. (1990). *Categorical Longitudinal Data: Loglinear Analysis of Panel, Trend and Cohort Data*. Newbury Park, CA: Sage.
- Heinen, T. (1996). *Latent Class and Discrete Latent Trait Models: Similarities and Differences*. Thousand Oaks, CA: Sage.
- Lanza, S. T., Collins, L. M., Lemmon, D. R., and Schafer, J. L. (2007). PROC LCA: A SAS procedure for latent class analysis. *Structural Equation Modeling* **14**, 671–694.
- Lazarsfeld, P. F. (1950). The logical and mathematical foundation of latent structure analysis and the interpretation and mathematical foundation of latent structure analysis. In Stouffer, S. A., Guttman, L., Suchman, E. A., *et al.* (eds.) *Measurement and Prediction*, pp 362–472. Princeton, NJ: Princeton University Press.
- Lindsay, B., Clogg, C. C., and Grego, J. (1991). Semiparametric estimation in the Rasch model and related models, including a simple latent class model for item analysis. *Journal of the American Statistical Association* **86**, 96–107.

- Magidson, J. and Vermunt, J. K. (2001). Latent class factor and cluster models, bi-plots and related graphical displays. *Sociological Methodology* **31**, 223–264.
- McLachlan, G. J. and Peel, D. (2000). *Finite Mixture Models*. New York: Wiley.
- Proctor, C. H. (1970). A probabilistic formulation and statistical analysis of Guttman scaling. *Psychometrika* **35**, 73–78.
- Rost, J. (1990). Rasch models in latent classes. An integration of two approaches to item analysis. *Applied Psychological Measurement* **14**, 271–282.
- Sijtsma, K. and Molenaar, I. W. (2002). *Introduction to Nonparametric Item Response theory*. Thousand Oaks, CA: Sage.
- Skrondal, A. and Rabe-Hesketh, S. (2004). *Generalized Latent Variable Modeling: Multilevel, Longitudinal and Structural Equation Models*. London: Chapman and Hall/CRC.
- Van de Pol, F. and Langeheine, R. (1990). Mixed Markov latent class models. *Sociological Methodology* **20**, 213–247.
- Vermunt, J. K. (1997). *Log-Linear Models for Event Histories*. Thousand Oaks, CA: Sage.
- Vermunt, J. K. (2001). The use restricted latent class models for defining and testing nonparametric and parametric IRT models. *Applied Psychological Measurement* **25**, 283–294.
- Vermunt, J. K. (2003). Multilevel latent class models. *Sociological Methodology* **33**, 213–239.
- Vermunt, J. K. (2007). Growth models for categorical response variables: Standard, latent-class, and hybrid approaches. In van Montfort, K., Oud, H., and Satorra, A. (eds.) *Longitudinal Models in the Behavioral and Related Sciences*, pp 139–158. Mahwah, NJ: Erlbaum.
- Vermunt, J. K. and Magidson, J. (2002). Latent class cluster analysis. In Hagenars, J. and McCutcheon, A. (eds.) *Applied Latent Class Analysis*, pp 89–106. Cambridge University Press.
- Vermunt, J. K. and Van Dijk, L. (2001). A nonparametric random-coefficients approach: The latent class regression model. *Multilevel Modelling Newsletter* **13**, 6–13.
- Wedel, M. and DeSarbo, W. (1994). A review of recent developments in latent class regression models. In Bagozzi, R. P. (ed.) *Advanced Methods of Marketing Research*, pp 352–388. Cambridge, MA: Blackwell.
- Wolfe, J. H. (1970). Pattern clustering by multivariate cluster analysis. *Multivariate Behavioral Research* **5**, 329–350.
- Yung, Y. F. (1997). Finite mixtures in confirmatory factor-analysis models. *Psychometrika* **62**, 297–330.

## Further Reading

- Hagenars, J. A. and McCutcheon, A. L. (2002). *Applied Latent Class Analysis*. Cambridge, UK: Cambridge University Press.
- Magidson, J. and Vermunt, J. K. (2004). Latent class models. In Kaplan, D. (ed.) *The Sage Handbook of Quantitative Methodology for the Social Sciences*, ch. 10, pp 175–198. Thousand Oakes, CA: Sage.
- Vermunt, J. K. (forthcoming). Mixture models for multilevel data sets. In Hox, J. and Roberts, J. K. (eds.) *The Handbook of Advanced Multilevel Analysis*.
- Vermunt, J. K., Tran, B., and Magidson, J. (2008). Latent class models in longitudinal research. In Menard, S. (ed.) *Handbook of Longitudinal Research: Design, Measurement, and Analysis*, pp 373–385. Burlington, MA: Elsevier.