



Model Identifiability

By *Guan-Hua Huang*

Keywords: *factor analysis, Fisher information matrix, latent class model, Markov chain Monte Carlo*

Abstract: A model is identifiable if there is a one-to-one correspondence between the probability distribution of the data and the values of model parameters. When applying a nonidentifiable model, different people may draw different conclusions from the same model of the observed data. Before one can meaningfully discuss the estimation of a model, model identifiability must be verified.

In some statistical models, different parameter values can give rise to identical probability distributions. When this happens, there will be a number of different parameter values associated with the maximum likelihood of any set of observed data. This is referred to as the *model identifiability problem*. For example, suppose someone attempts to compute the regression equation predicting Y from three variables X_1 , X_2 , and their sum ($X_1 + X_2$), the program will probably crash or give an error message because it cannot find a unique solution. The model is the same if $Y = 0.5X_1 + 1.0X_2 + 1.5(X_1 + X_2)$, $Y = 1.0X_1 + 1.5X_2 + 1.0(X_1 + X_2)$, or $Y = 2.0X_1 + 2.5X_2 + 0.0(X_1 + X_2)$; indeed, there are an infinite number of equally good possible solutions. Model identifiability is a particular problem for the latent class model, a statistical method for finding the underlying traits from a set of psychological tests because, by postulating latent variables, it is easy to introduce more parameters into a model than can be fitted from the data.

A model is identifiable if the parameter values uniquely determine the probability distribution of the data and the probability distribution of the data uniquely determines the parameter values. Formally, let ϕ be the parameter value of the model, y be the observed data, and $F(y; \phi)$ be the probability distribution of the data. A model is identifiable if for all $(\phi_0, \phi) \in \Phi$ and all $y \in S_Y$:

$$F(y; \phi_0) = F(y; \phi) \text{ if and only if } \phi_0 = \phi \quad (1)$$

where Φ denotes the set of all possible parameter values and S_Y the set of all possible values of the data.

The most common cause of model nonidentifiability is a poorly specified model. If the number of unique model parameters exceeds the number of independent pieces of observed information, the model is not identifiable. Consider the example of a latent class model that classifies people into three states (severely depressed/mildly depressed/not depressed) and that is used to account for the responses of a group of people to three psychological tests with binary (positive/negative) outcomes. Let (Y_1, Y_2, Y_3) denote the test results and let each take the value 1 when the outcome is positive and 0 when it is negative. S specifies the unobservable states where $S = 1$ where there is no depression, 2 where the depression is mild, and

National Chiao Tung University, Hsinchu, Taiwan

Update based on original article by Guan-Hua Huang, Wiley StatsRef: Statistics Reference Online ©2014 John Wiley & Sons, Ltd.

Wiley StatsRef: Statistics Reference Online, © 2014–2016 John Wiley & Sons, Ltd.
This article is © 2016 John Wiley & Sons, Ltd.
DOI: 10.1002/9781118445112.stat06411.pub2





3 where the depression is severe. The probability of the test results is then

$$\Pr(Y_1 = y_1, Y_2 = y_2, Y_3 = y_3) = \sum_{j=1}^3 \Pr(S = j) \prod_{m=1}^3 \Pr(Y_m = 1|S = j)^{y_m} \times \Pr(Y_m = 0|S = j)^{1-y_m} \quad (2)$$

The test results have $2^3 - 1 = 7$ independent patterns and the model requires 11 unique parameters (two probabilities for depression status $\Pr(S = 3)$, $\Pr(S = 2)$, and one conditional probability $\Pr(Y_m = 1|S = j)$ for each depression status j and test m); therefore, the model is not identifiable.

If the model is not identifiable, one can make it so by imposing various constraints upon the parameters. When there appears to be sufficient total observed information for the number of estimated parameters, it is also necessary to specify the model unambiguously. For the above-mentioned latent class model, suppose that, for the second and third tests, the probabilities of observing a positive test result are the same for people with severe, mild, or no depression (i.e., $\Pr(Y_m = 1|S = 3) = \Pr(Y_m = 1|S = 2) = \Pr(Y_m = 1|S = 1) = p_m$ for $m = 2, 3$). In other words, only the first test discriminates between the unobservable states of depression. The model now has only seven parameters, which is equal to the number of independent test result patterns. The probability distribution of test results becomes

$$\Pr(Y_1 = y_1, Y_2 = y_2, Y_3 = y_3) = \Theta \prod_{m=2}^3 (p_m)^{y_m} (1 - p_m)^{1-y_m} \quad (3)$$

where

$$\Theta = (1 - \eta_2 - \eta_3)(p_{11})^{y_1} (1 - p_{11})^{1-y_1} + \eta_2(p_{12})^{y_1} (1 - p_{12})^{1-y_1} + \eta_3(p_{13})^{y_1} (1 - p_{13})^{1-y_1} \quad (4)$$

$\eta_2 = \Pr(S = 2)$, $\eta_3 = \Pr(S = 3)$, $p_{11} = \Pr(Y_1 = 1|S = 1)$, $p_{12} = \Pr(Y_1 = 1|S = 2)$, and $p_{13} = \Pr(Y_1 = 1|S = 3)$. Θ imposes two restrictions on parameters (i.e., for $y_1 = 1$ or 0), and there are five parameters to consider (i.e., η_2 , η_3 , p_{11} , p_{12} , and p_{13}). Because the number of restrictions is less than the number of parameters of interest, Θ and the above-mentioned latent class model are not identifiable—the same probability distributions could be generated by supposing that there was a large chance of being in a state with a small effect on the probability of being positive on test 1 or by supposing that there was a small chance of being in this state but it was associated with a large probability of responding positively.

Sometimes, it is difficult to find an identifiable model. A weaker form of identification, called *local identifiability*, may exist, namely, it may be that other parameters generate the same probability distribution as ϕ_0 does, but one can find an open neighborhood of ϕ_0 that contains none of these parameters^[1]. For example, we are interested in β in the regression $Y = \beta^2 X$ (the square root of the association between Y and X). $\beta = 1$ and $\beta = -1$ result in the same Y prediction; thus, the model is not (globally) identifiable. However, the model is locally identifiable because one can easily find two nonoverlapping intervals $(0.5, 1.5)$ and $(-1.5, -0.5)$ for 1 and -1 , respectively. Specifically, a distribution F is locally identifiable at the parameter ϕ_0 if there exists some neighborhood χ of ϕ_0 , such that for all $\phi \in \chi \subset \Phi$ and all \mathbf{y}

$$F(\mathbf{y}; \phi) = F(\mathbf{y}; \phi_0) \text{ if and only if } \phi_0 = \phi \quad (5)$$

A locally but not globally identifiable model does not have a unique interpretation, but one can be sure that, in the neighborhood of the selected solution, there exist no other equally good solutions; thus, the problem is reduced to determining the regions where local identifiability applies. This concept is especially useful in models containing nonlinearities as the above-mentioned regression example, or models with complex structures, for example, *factor analysis*^[2] and regression extension of latent class models^[3].





It is difficult to specify general conditions that are sufficient to guarantee (global) identifiability. Fortunately, it is fairly easy to determine local identifiability. One can require that the columns of the Jacobian matrix, the first-order partial derivative of the likelihood function with respect to the unique model parameters, are independent^[1,4]. Alternatively, we can examine whether the *Fisher information matrix* possesses *eigenvalues* > 0 ^[5]. Formann^[6] showed that these two approaches are equivalent. A standard practice for checking local identifiability involves using multiple sets of initial values for parameter estimation. Different sets of initial values that yield the same likelihood maximum should result in the same final parameter estimates. If not, the model is not locally identifiable.

Complications often arise from applying the above-mentioned methods to a given analysis. Ideally, one would want to determine those regions of the parameter space in which a given model is locally identifiable. Because this is typically computationally difficult, these methods are often evaluated with respect to estimated parameters to establish local identifiability at estimated values^[4]. When using the *Fisher information matrix*, there is one more complication. As the observed Fisher information matrix, the negative matrix of the second-order partial derivatives of the log likelihood, is typically used to estimate the standard errors of maximum likelihood estimators^[7], the *Fisher information matrix* is not always obtained and the observed Fisher information matrix is used for empirical checking. Empirical identifiability checking through the observed Fisher information might cause errors because we use the *single* observation in place of the averaged effect. It needs to be implemented cautiously.

When applying a nonidentifiable model, different people may draw different conclusions from the same model of the observed data. Before one can meaningfully discuss the estimation of a model, model identifiability must be verified. If researchers come up against identifiability problems, they can first identify the parameters involved in the lack of identifiability from their extremely large asymptotic standard errors^[6] and then impose reasonable constraints on identified parameters based on prior knowledge or empirical information.

In the past decades, the development of Markov chain Monte Carlo (MCMC) methods and progress of computer technology facilitate the popularity of performing Bayesian analysis for latent class models. In Bayesian setting, there are two sources of nonidentifiability, the typical parameter identifiability problem as discussed earlier and the well-known likelihood invariance under label switching. The parameter nonidentifiability causes no real difficulties for Bayesian analysis as it is always possible to resolve the issue via placing proper prior distributions on the unidentified parameters^[8]. Regarding the label switching problem, notice that the probability (Equation 2) is invariant when switching state labels. Namely, if the values of (conditional) probabilities are exchanged between $S = 1$ and $S = 2$, the probability (Equation 2) remains the same, which leads to nonidentifiability of the labels of the state variable S . This is not a problem for a deterministic algorithm such as the maximum likelihood or EM (expectation–maximization), but it complicates the inference from sampling procedures such as MCMC because the labels of states may be randomly switched during the iterative process. If the prior distributions for model parameters do not distinguish the states of Equation (2), the resulting posterior distributions will be invariant to all permutations of state labels. Hence, the ergodic averages over the MCMC samples from the posterior distributions are meaningless. Many approaches have been proposed to deal with the label switching problem in Bayesian latent class analysis. The most commonly used approach is to impose some artificial ordering constraints on model parameters. In addition, researchers develop various algorithms that relabel the state statuses in each MCMC iteration using k -means-type approaches or label-invariant loss functions^[9].

Related Articles

Latent Class Analysis; Bayesian Inference; Markov Chain Monte Carlo Algorithms; Identifiability; Markov Chain Monte Carlo (MCMC).



References

- [1] McHugh, R.B. (1956) Efficient estimation and local identification in latent class analysis. *Psychometrika* **21**, 331–347.
- [2] Shapiro, A. (1985) Identifiability of factor analysis: some results and open problems. *Linear Algebra Appl.* **70**, 1–7.
- [3] Huang, G.H. and Bandeen-Roche, K. (2004) Building an identifiable latent variable model with covariate effects on underlying and measured variables. *Psychometrika* **69**, 5–32.
- [4] Goodman, L.A. (1974) Exploratory latent structure analysis using both identifiable and unidentifiable models. *Biometrika* **61**, 215–231.
- [5] Rothenberg, T.J. (1971) Identification in parametric models. *Econometrica* **39**, 577–591.
- [6] Formann, A.K. (1992) Linear logistic latent class analysis for polytomous data. *J. Am. Stat. Assoc.* **87**, 476–486.
- [7] Efron, B. and Hinkley, D.V. (1978) Assessing the accuracy of the maximum likelihood estimator: observed versus expected Fisher information. *Biometrika* **65**, 457–487.
- [8] Lindley, D.V. (1971) *Bayesian Statistics: A Review*, SIAM, Philadelphia, PA.
- [9] Jasra, A., Holmes, C.C., and Stephens, D.A. (2005) Markov chain Monte Carlo methods and the label switching problem in Bayesian mixture modeling. *Stat. Sci.* **20**, 50–67.