CHAPTER 1. DISCRETE RESPONSE MODELS

1. INTRODUCTION

When economic behavior is expressed as a continuous variable, a linear regression model is often adequate to describe the impact of economic factors on this behavior, or to predict this behavior in altered circumstances. For example, a study of food expenditures as a function of price indices for commodity groups and income, using households from the Consumer Expenditure Survey, can start by modeling indirect utility as a translog function and from this derive a linear in logs regression equation for food expenditures that does a good job of describing behavior. This situation remains true even when the behavioral response is limited in range (e.g., food consumption of households is non-negative) or integer-valued (e.g., college enrollment by state), provided these departures from a unrestricted continuous variable are not conspicuous in the data (e.g., food consumption is observed over a range where the non-negativity restriction is clearly not binding; college enrollments are in the thousands, so that round-off of the dependent variable to an integer is negligible relative to other random elements in the model). However, there are a variety of economic behaviors where the continuous approximation is not a good one. Here are some examples:

- (1) For individuals: Whether to attend college; whether to marry; choice of occupation; number of children; whether to buy a house; what brand of automobile to purchase; whether to migrate, and if so where; where to go on vacation.
- (2) For firms: Whether to build a plant, and if so, at what location; what commodities to produce; whether to shut down, merge or acquire other firms; whether to go public or private; whether to accept union demands or take a strike.

For sound econometric analysis, one needs probability models that approximate the true data generation process. To find these, it is necessary to think carefully about the economic behavior, and about the places where random factors enter this behavior. For simplicity, we initially concentrate on a single binomial (Yes/No) response. An example illustrates the process:

Yellowstone National Park has been overcrowded in recent years, and large user fees to control demand are under consideration. The National Park Service would like to know the elasticity of demand with respect to user fees, and the impact of a specified fee increase on the total number of visitors and on the visitors by income bracket. The results of a large household survey are available giving household characteristics (income, number of children, etc.), choice of vacation site, and times and costs associated with vacations at alternative sites. Each vacation is treated as an observation.

Start with the assumption that households are utility maximizers. Then, each household will have an indirect utility function, *conditioned* on vacation site, that gives the payoff to choosing this

particular site and then optimizing consumption in light of this choice. This indirect utility function will depend on commodity prices and on household income net of expenditures mandated by the vacation site choice. It may also contain factors such as household tastes and perceptions, and unmeasured attributes of sites, that are, from the standpoint of the analyst, random. (Some of what appears to be random to the analyst may just be heterogeneity in tastes and perceptions over the population.) Now consider the difference between the indirect utility of a Yellowstone vacation and the *maximum* indirect utilities of alternative uses of leisure. This is a function $y^* = f(z,\zeta)$ of observed variables z and unobserved variables ζ. We put a "*" on the utility difference y to indicate that is *latent* rather than observed directly. Included in z are variables such as household income, wage rate, family characteristics, travel time and cost to Yellowstone, and so forth. The form of this function will be governed by the nature of indirect utility functions and the sources of ζ . In some applications. it makes sense to parameterize the initial indirect utility functions tightly, and then take f to be the function implied by this. Often, it is more convenient to take f to be a form that is flexibly parameterized and convenient for analysis, subject only to the generic properties that a difference of indirect utility functions should have. In particular, it is almost always possible to approximate f closely by a function that is linear in parameters, with an additive disturbance: $f(z,\zeta) \approx x\beta - \varepsilon$, where β is a k×1 vector of unknown parameters, x is a 1×k vector of transformations of z, and $\varepsilon = -f(z,\zeta) + \mathbf{E}f(z,\zeta)$ is the deviation of f from its expected value in the population. Such an approximation might come, for example, from a Taylor's expansion of $\mathbf{E}f$ in powers of (transformed) observed variables z.

Suppose the gain in utility from vacationing in Yellowstone rather than at an alternative site is indeed given by $y^* = x\beta - \epsilon$. Suppose the disturbance ϵ is known to the household and unknown to the econometrician, but the cumulative distribution function (CDF) of ϵ is a function $F(\epsilon)$ that is known up to a finite parameter vector. The utility-maximizing household will then choose Yellowstone if $y^* > 0$, or $\epsilon < x\beta$. The probability that this occurs, given x, is

$$P(\varepsilon < x\beta) = F(x\beta).$$

Define y = 1 if Yellowstone is chosen, y = -1 otherwise; then, y is an (observed) indicator for the event $y^* > 0$. The probability law governing observed behavior is then, in summary,

$$P(y|x\beta) = \begin{cases} F(x\beta) & if y = 1 \\ 1 - F(x\beta) & if y = -1 \end{cases}.$$

Assume that the distribution of ε is symmetric about zero, so that $F(\varepsilon) = 1 - F(-\varepsilon)$; this is not essential, but it simplifies notation. The probability law then has an even more compact form,

$$P(y|x\beta) = F(yx\beta)$$
.

How can you estimate the parameters β ? An obvious approach is maximum likelihood. The log likelihood of an observation is

$$l(\beta | y,x) = \log P(y|x\beta) = \log F(yx\beta)$$
.

If you have a random sample with observations t = 1,...,T, then the sample log likelihood is

$$L_{T}(\beta) = \sum_{t=1}^{T} \log F(y_{t}x_{t}\beta).$$

The associated score and hessian of the log likelihood are

$$\nabla_{\beta} L_{T}(\beta) = \sum_{t=1}^{T} y_{t} x_{t}' F'(y_{t} x_{t} \beta) / F(y_{t} x_{t} \beta)$$

$$\nabla_{\beta\beta} L_{T}(\beta) = \sum_{t=1}^{T} x_{t}' x_{t} \{F''(y_{t}x_{t}\beta)/F(y_{t}x_{t}\beta) - [F'(y_{t}x_{t}\beta)/F(y_{t}x_{t}\beta)]^{2} \}.$$

A maximum likelihood program will either ask you to provide these formula, or will calculate them for you analytically or numerically. If the program converges, then it will then find a value of β (and any additional parameters upon which F depends) that are (at least) a local maximum of L_T . It can fail to converge to a maximum if no maximum exists or if there are numerical problems in the evaluation of expressions or in the iterative optimization. The estimates obtained at convergence will have the usual large-sample properties of MLE, provided the usual regularity conditions are met, as discussed later.

It is sometimes useful to write the score and hessian in a slightly different way. Let d = (y+1)/2; then d = 1 for Yellowstone, d = 0 otherwise, and d is an indicator for a Yellowstone trip. Then, we can write

$$l(y|x,\beta) = d \cdot \log F(x\beta) + (1-d) \cdot \log F(-x\beta).$$

Differentiating this expression, and noting that $F'(x\beta) = F'(-x\beta)$, we get

$$\nabla_{\beta} l = xF'(x\beta) \{d/F(x\beta) - (1-d)/F(-x\beta)\} = w(x\beta) \cdot x \cdot [d - F(x\beta)],$$

where $w(x\beta) = F'(x\beta)/F(x\beta)F(-x\beta)$. The sample score is then

$$\nabla_{\beta} L_{T}(\beta) = \sum_{t=1}^{T} w(x_{t}\beta) \cdot x_{t}' \cdot [d_{t} - F(x_{t}\beta)].$$

The MLE condition that the sample score equal zero can be interpreted as a weighted *orthogonality* condition between a residual $[d - F(x\beta)]$ and the explanatory variables x. Put another way, a weighted non-linear least squares (NLLS) regression $d_t = F(x_t\beta) + \eta_t$, with observation t weighted by $w(x_t\beta)^{1/2}$, will be equivalent to MLE.¹

The hessian can also be rewritten using d rather than y: $\nabla_{\beta\beta}l = -x'x \cdot s(x\beta)$, where $s(x\beta) =$

$$\frac{F'(x\beta)^2}{F(x\beta)F(-x\beta)} - [\mathsf{d} - \mathsf{F}(x\beta)] \left\{ \frac{F''(x\beta)}{F(x\beta)F(-x\beta)} - \frac{F'(x\beta)^2(1-2F(x\beta))}{F(x\beta)^2F(-x\beta)^2} \right\}.$$
 The expectation of $\mathsf{s}(\mathsf{x}\beta)$ at

¹To be precise, iterated NLLS, with the β appearing in the weighting function replaced by the last iterate, will converge to the MLE estimator; a single NLLS *without weighting* provides estimates of β that are consistent and asymptotically normal, but not asymptotically efficient; and *one* iterate with weights calculated from a consistent estimator of β will be asymptotically equivalent to MLE.

the true value
$$\beta_0$$
 is $\frac{F'(x\beta_o)^2}{F(x\beta_o)F(-x\beta_o)} > 0$, so that the sample sum of the hessians of the observations

in sufficiently large samples is eventually almost surely negative definite in a neighborhood of β_o . It should be clear from the sample score, or the analogous NLLS regression, that the distribution function F enters the likelihood function in an intrinsic way. Unlike linear regression, there is no simple estimator of β that rests only on assumptions about the first two moments of the disturbance distribution.²

2. FUNCTIONAL FORMS AND ESTIMATORS

In principle, the CDF $F(\epsilon)$ will have a form deduced from the application; in many cases, this form would naturally be conditioned on the observed explanatory variables. However, an almost universal practice is to assume that $F(\epsilon)$ has one of the following standard distributions that are not conditioned on x:

- (1) *Probit*: F is standard normal.
- (2) *Logit*: $F(\varepsilon) = 1/(1+e^{-\varepsilon})$, the standard logistic CDF.
- (3) *Linear*: $F(\varepsilon) = \varepsilon$, for $0 \le \varepsilon \le 1$, the standard uniform distribution.
- (4) Log-Linear: $F(\varepsilon) = e^{\varepsilon}$, for $\varepsilon \le 0$, a standard exponential CDF.

There are many canned computer programs to fit models (1) or (2). Model (3) can be fit by linear regression, although heteroscedasticity is an issue. Model (4) is not usually a canned program when one is dealing with individual observations, but for repeated observations at each configuration of x it is a special case of the *discrete analysis of variance* model that is widely used in biostatistics and can be fitted using ANOVA or regression methods. Each of the distributions above has the property that the function $s(x\beta)$ that appears in the hessian is globally positive, so that the log likelihood function is globally concave. This is convenient in that any local maximum is the global maximum, and any stable hill-climbing algorithm will always get to the global maximum. The linear and log-linear distributions are limited in range. This is typically not a problem if the range of x is such that the probabilities are bounded well away from zero and one, but can be a serious problem when some probabilities are near or at the extremes, particularly when the model is used for forecasting.

The remainder of this section deals with some alternatives to maximum likelihood estimation, and can be skipped on first reading. Recall that MLE chooses the parameter vector β to achieve orthogonality between the explanatory variables x, and residuals d - $F(x\beta)$, with weights $w(x\beta)$. When the explanatory variables are grouped, or for other reasons there are multiple responses observed for the same x, there is another estimation procedure that is useful. Let j = 1,...,J index the possible x configurations, m_i denote the number of responses observed at configuration x_i , and s_i

²We will see later that there are some more robust estimators, not as simple, that avoid having to place F in a parametric family, or use a non-parametric estimate of F. Sometimes assumptions on F are sufficiently problematic so this extra complexity is worth the trouble.

denote the number of "successes" among these responses (i.e., the number with d=1). Let $p_j=F(x_j\beta_o)$ denote the true probability of a success at configuration x_j . Invert the CDF to obtain $c_j=F^{-1}(p_j)=x_j\beta$. Note that p=F(c) implies $\partial c/\partial p=1/F'(c)$ and $\partial^2 c/\partial p^2=-F''(c)/F'(c)^3$. Then, a Taylor's expansion of $F^{-1}(s_i/m_i)$ about p_i gives

$$F^{-1}(s_j/m_j) = F^{-1}(p_j) + \frac{s_j/m_j - p_j}{F'(F^{-1}(p_j))} - \frac{(s_j/m_j - p_j)^2}{2} \cdot \frac{F''(F^{-1}(q_j))}{F'(F^{-1}(q_i))^3} = x_j\beta + v_j + \xi_j,$$

where q_j is a point between p_j and s_j/m_j , $v_j = (s_j/m_j - p_j)/F'(F^{-1}(p_j))$ is a disturbance that has expectation zero and a variance proportional to $p_j(1-p_j)/m_j$, and ξ_j is a disturbance that goes to zero in probability relative to v_j . Then, when the m_j are all large (the rule-of-thumb is $s_j \ge 5$ and $m_j - s_j \ge 5$), the regression

$$F^{-1}(s_i/m_i) = x_i\beta + v_i$$

gives consistent estimates of β . This is called *Berkson's method*. It can be made asymptotically equivalent to MLE if a FGLS transformation for heteroscedasticity is made. Note however that in general this transformation is not even defined unless s_j is bounded away from zero and m_j , so it does not work well when some x's are continuous and cell counts are small. Note that Berkson's transformation in the case of probit is $\Phi^{-1}(s_j/m_j)$; in the case of logit is $\log(s_j/(m_j-s_j))$; in the case of linear is s_j ; and in the case of the exponential model is $\log(s_j/m_j)$. It is a fairly general proposition that the asymptotic approximation is improved by using the transformation $F^{-1}((s_j+0.5)/(m_j+1))$ rather than $F^{-1}(s_j/m_j)$ as the dependent variable in the regression; for logit, this minimizes the variance of the second-order error.

There is an interesting connection between the logit model and a technique called *normal linear discriminant analysis*. Suppose that the conditional distributions of x, given d=1 or given d=0, are both multivariate normal with respective mean vectors μ_1 and μ_0 , and a *common* covariance matrix Ω . Note that these assumptions are not necessarily very plausible, certainly not if some of the x variables are limited or discrete. If the assumptions hold, then the means μ_0 and μ_1 and the covariance matrix Ω can be estimated from sample averages, and by Bayes law the conditional distribution of d given x when a proportion q_1 of the population has state d=1 has a logit form

$$P(d=1|x) = \frac{q_1 n(x - \mu_1, \Omega)}{q_0 n(x - \mu_0, \Omega) + q_1 n(x - \mu_1, \Omega)} = \frac{1}{1 + \exp(-\alpha - x\beta)},$$

where $\beta = \Omega^{-1}(\mu_1 - \mu_0)$ and $\alpha = \mu_1' \Omega^{-1} \mu_1 - \mu_0' \Omega^{-1} \mu_0 + \log(q_1/q_0)$. This approach produces a fairly robust (although perhaps inconsistent) estimator of the logit parameters, even when the normality assumptions are obviously wrong.

3. STATISTICAL PROPERTIES OF MLE

The MLE estimator for most binomial response models is a special case of the general setup treated in the statistical theory of MLE, so that the incantation "consistent and asymptotically normal (CAN) under standard regularity conditions" is true. This is a simple enough application so that it

is fairly straightforward to see what these "regularity" conditions mean, and verify that they are satisfied. This is a thought exercise worth going through whenever you are applying the maximum likelihood method. First, here is a list of fairly general sufficient conditions for MLE to be CAN in discrete response models; these are taken from McFadden "Quantal Response Models", <u>Handbook of Econometrics</u>, Vol. 2, p. 1407. Commentaries on the assumptions are given in italics.

- (1) The domain of the explanatory variables is a measurable set X with a probability p(x). This just means that the explanatory variables have a well-defined distribution. It certainly holds if the domain (support) of X is a closed set, and p is a continuous density on X.
- (2) The parameter space is a subset of \mathbb{R}^k , and the true parameter vector is in the interior of this space. This says you have a finite-dimensional parametric problem. This assumption does not require that the parameter space be bounded, in contrast to many sets of assumptions used to conclude that MLE are CAN. The restriction that the true parameter vector be in the interior excludes some cases where CAN breaks down. This is not a restrictive assumption in most applications, but it is for some. For example, suppose a parameter in the probit model is restricted (by economic theory) to be non-negative, and that this parameter is in truth zero. Then, its asymptotic distribution will be the (non-normal) mixture of a half-normal and a point mass.
- (3) The response model is measurable in x, and for almost all x is continuous in the parameters. The standard models such as probit, logit, and the linear probability model are all continuous in their argument and in x, so that the assumption holds. Only pathological applications in which a parameter determines a "trigger level" will violate this assumption.
- (4) The model satisfies a global identification condition (that guarantees that there is at most one global maximum; see McFadden, *ibid*, p. 1407). The concavity of the log likelihood of an observation for probit, logit, linear, and log linear models guarantees global identification, provided only that the x's are not linearly dependent.
- (5) The model is once differentiable in the parameters in some neighborhood of the true values. This is satisfied by the four CDF from Section 2 (provided parameters do not give observations on the boundary in the linear or log linear models where probabilities are zero or one), and by most applications. This is weaker than most general MLE theorems, which assume the log likelihood is twice or three times continuously differentiable.
- (6) The log likelihood and its derivative have bounds independent of the parameters in some neighborhood of the true parameter values. The first derivative has a Lipschitz property in this neighborhood. *This property is satisfied by the four CDF, and any CDF that are continuously differentiable.*

(7) The information matrix, equal to the expectation of the outer product of the score of an observation, is nonsingular at the true parameters. *This is satisfied automatically by the four CDF in Section 2, provided the* x's are not linearly dependent.

The result that conditions (1)-(7) guarantee that MLE estimates of β are CAN is carried out essentially by linearizing the first-order condition for the estimator using a Taylor's expansion, and arguing that higher-order terms than the linear term are asymptotically negligible. With lots of differentiability and uniform bounds, this is an easy argument. A few extra tricks are needed to carry this argument through under the weaker smoothness conditions contained in (1)-(7).

4. EXTENSIONS OF THE MAXIMUM LIKELIHOOD PRINCIPLE

The assumptions under which the maximum likelihood criterion produces CAN estimates include, critically, the condition (2) that the parametric family of likelihoods that are being maximized include the true data generation process. There are several reasons that this assumption can fail. First, you may have been mistaken in your assumption that the model you have written down includes the truth. This might happen in regression analysis because some variable that you think does not influence the dependent variable or is uncorrelated with the included variables actually does belong in the regression. Or, in modeling a binomial discrete response, you may assume that the disturbance in the model $y^* = x\beta - \epsilon$ is standard normal when it is in truth logistic. Second, you may deliberately write down a model you suspect is incorrect, simply because it is convenient for computation or reduces data collection problems. For example, you might write down a model that assumes observations are independent even though you suspect they are not. This might happen in discrete response analysis where you observe several responses from each economic agent, and suspect there are unobserved factors such as tastes that influence all the responses of this agent.

What are the statistical consequences of this model misspecification? The answer is that this will generally cause the CAN property to fail, but in some cases the failure is less disastrous than one might think. The most benign situation arises when you write down a likelihood function that fails to use all the available data in the most efficient way, but is otherwise consistent with the true likelihood function. For example, if you have several dependent variables, such as binomial responses on different dates, you may write down a model that correctly characterizes the marginal likelihood of each response, but fails to characterize the dependence between the responses. This setup is called *quasi-maximum likelihood* estimation. What may happen in this situation is that not all the parameters in the model will be identified, but those that are identified are estimated CAN, although not necessarily with maximum efficiency. In the example, it will be parameters characterizing the correlations across responses that are not identified. Also fairly benign is a method called *pseudo-maximum likelihood* estimation, where you write down a likelihood function with the property that the resulting maximum likelihood estimates are in fact functions only of selected moments of the data. A classic example is the normal regression model, where the maximum likelihood estimates depend only on first and second moments of the data. Then the estimates that come out of this criterion will be CAN even if the pseudo-likelihood function is misspecified, so long as the true likelihood function and the pseudo-likelihood function coincide for the moments that the estimators actually use.

More tricky is the situation where the likelihood you write down is not consistent with the true likelihood function. In this case, the parameters in the model you estimate will not necessarily match up, even in dimension, with the parameters of the true model, and there is no real hope that you will get reasonable estimates of these true parameters. However, even here there is an interesting result. Under quite general conditions, it is possible to talk about the "asymptotically least misspecified model", defined as the model in your misspecified family that asymptotically has the highest log likelihood. To set notation, suppose f(y|x) is the true data generation process, and $g(y|x,\beta)$ is the family of misspecified models you consider. Define β_1 to be the parameters that maximize

$$\mathbf{E}_{\mathbf{y},\mathbf{x}} f(\mathbf{y}|\mathbf{x}) \cdot \log g(\mathbf{y}|\mathbf{x},\beta).$$

Then, β_1 determines the least misspecified model. While β_1 does not characterize the true data generation process, and the parameters as such may even be misleading in describing this process, what is true is that β_1 characterizes the model g that in a "likelihood metric" is as close an approximation as one can reach to the true data generation process when one restricts the analysis to the g family. Now, what is interesting is that the maximum likelihood estimates b from the misspecified model are CAN for β_1 under mild regularity conditions. A colloquial way of putting this is that MLE estimates are usually CAN for whatever it is they converge to in probability, even if the likelihood function is misspecified.

All of the estimation procedures just described, quasi-likelihood maximization, pseudo-likelihood maximization, and maximization of a misspecified likelihood function, can be interpreted as special cases of a general class of estimators called *generalized method of moment estimators*. One of the important features of these estimators is that they have asymptotic covariance matrices of the form $\Gamma^{-1}\Sigma\Gamma'^{-1}$, where Γ comes from the hessian of the criterion function, and Σ comes from the expectation of the outer product of the gradient of the criterion function. For true maximum likelihood estimation, this form reduces to Σ^{-1} , but more generally the full form $\Gamma^{-1}\Sigma\Gamma'^{-1}$ is required.

One important family of quasi-maximum likelihood estimators arises when an application has a likelihood function in two sub-vectors of parameters, and it is convenient to obtain preliminary CAN estimates of one sub-vector, perhaps by maximizing a conditional likelihood function. Then, the likelihood is maximized in the second sub-vector of parameters after plugging in the preliminary estimates of the first sub-vector. This will be a CAN procedure under general conditions, but it is necessary to use a formula of the form $\Gamma^{-1}\Sigma\Gamma'^{-1}$ for its asymptotic covariance matrix, where Σ includes a contribution from the variance in the preliminary estimates of the first sub-vector. The exact formulas and estimators for the terms in the covariance matrix are given in the lecture notes on generalized method of moments.

5. TESTING HYPOTHESES

It is useful to see how the general theory of large sample hypothesis testing plays out in the discrete response application. For motivation, return to the example of travel to Yellowstone Park. The basic model might be binomial logit,

$$P(y|x\beta) = F(yx\beta) = 1/(1 + \exp(-yx\beta)),$$

where x includes travel time and travel cost to Yellowstone, and family income, all appearing linearly:

$$x\beta = TT \cdot \beta_1 + TC \cdot \beta_2 + I \cdot \beta_3 + \beta_4$$

with TT = travel time, TC = travel cost, I = income. The parameter β_4 is an intercept term that captures the "average" desirability of Yellowstone relative to alternatives after travel factors have been taken into account. The Park Service is particularly concerned that an increase in Park entry fees, which would increase overall travel cost, will have a particularly adverse effect on low income families, and asks you to test the hypothesis that sensitivity to travel cost increases as income falls. This suggests the alternative model

$$x\beta = TT \cdot \beta_1 + TC \cdot \beta_2 + I \cdot \beta_3 + \beta_4 + \beta_5 \cdot TC/I$$

with the null hypothesis that $\beta_5 = 0$. This hypothesis can be tested by estimating the model without the null hypothesis imposed, so that β_5 is estimated. The Wald test statistic is the quadratic form (b_5 $-0)'V(b_5)^{-1}(b_5-0)$; it is just the square of the T-statistic for this one-dimensional hypothesis, and it is asymptotically chi-square distributed with one degree of freedom when the null hypothesis is true. When the null hypothesis is non-linear or of higher dimension, the Wald statistic requires retrieving the covariance matrix of the unrestricted estimators, and forming the matrix of derivatives of the constraint functions evaluated at b. An alternative that is computationally easier when both the unrestricted and restricted models are easy to estimate is to form the Likelihood Ratio statistic $2[L_T(b) - L_T(b^*)]$, where b and b^* are the estimates obtained without the null hypothesis and with the null hypothesis imposed, respectively, and L_T is the sample log likelihood. This statistic is asymptotically equivalent to the Wald statistic. Finally, the Lagrange Multiplier statistic is obtained by estimating the model under the null hypothesis, evaluating the score of the unrestricted model at the restricted estimates, and then testing whether this score is zero. In our example, there is a slick way to do this. Regress a normalized residual $[d_t - F(x_t b)]/[F(x b)F(-x b)]^{1/2}$ from the restricted model on the weighted explanatory variables $x \cdot F'(xb)/[F(x b)F(-x b)]^{1/2}$. that appear in the unrestricted model. The F-test for the significance of the explanatory variables in this regression is asymptotically equivalent to the Lagrange Multiplier test. The reason this trick works is that the Lagrange Multiplier test is a test of orthogonality between the normalized residual and the weighted variables in the unrestricted model.

6. MULTINOMIAL RESPONSE

Conceptually, it is straightforward to move from modeling binomial response to modeling multinomial response. When consumers or firms choose among multiple, mutually exclusive alternatives, such as choice of brand of automobile, occupation, or plant location, it is natural to introduce the economic agent's objective function (utility for consumers, profit for firms), and assume that choice maximizes this objective function. Factors unobserved by the analyst, particularly heterogeneity in tastes or opportunities, can be interpreted as random components in the

objective functions, and choice probabilities derived as the probabilities that these unobserved factors are configured so as to make the respective alternatives optimal.

Suppose there are J alternatives, indexed $C = \{1,...,J\}$, and suppose the economic agent seeks to maximize an objective function $U(z_i,s,v_i)$, where z_i are observed attributes of alternative i, s are characteristics of the decision maker, and v_i summarizes all the unobserved factors that influence the attractiveness of alternative i. Then, the multinomial response probability is

$$P_{C}(i|z,s) = Prob(\{v|U(z_{i},s,v_{i}) > U(z_{i},s,v_{i}) \text{ for } j \neq i\}),$$

where $z = (z_1,...,z_J)$. For example, if $C = \{1,...,J\}$ is the set of automobile brands, with z_i the attributes of brand i including price, size, horsepower, fuel efficiency, etc., then this model can be used to explain brand choice, or to predict the shares of brands as the result of changing prices or new model introductions. If one of the alternatives in C is the "no purchase" alternative, the model can describe the demand for cars as well as brand choice. If C includes both new and used alternatives, then it can explain replacement behavior. If $i \in C$ identifies a portfolio of two brands, or one brand plus a "no purchase", it can explain the holdings of two-car families.

Placing U in a parametric family and making v a random vector with a parametric probability distribution produces a parametric probability law for the observations. However, it is difficult to do this in a way that leads to simple algebraic forms that do not require multivariate integration. Consequently, the development of multinomial response models has tended to be controlled by computational issues, which may not accommodate some features that might seem sensible given the economic application, such as correlation of unobservables across alternative portfolios that have common elements.

The simplest multinomial response model is *multinomial logit* (MNL), which has a closed form

$$P_C(i|z,s) = \exp(x_i\beta) / \sum_{j \in C} \exp(x_j\beta),$$

where x_i is a vector of known functions of z_i and s. This model is derived from the maximizing framework above by assuming $U(z_i, s, v_i) = x_i \beta + \varepsilon_i$, with the ε_i independently identically distributed

with the special CDF $\exp(-e^{-\varepsilon_i})$, termed the *Type I extreme value distribution*.

The *likelihood* of observation n from a MNL model for choice from C is

$$l_n = \sum_{i \in C} d_{in} \cdot \log(P_{Cn}(i)),$$

where $P_{Cn}(i) = \exp(x_{in}\beta) / \sum_{k \in C} \exp(x_{kn}\beta)$, and $d_{in} = 1$ indicates choice and $d_{jn} = 0$ for non-chosen

alternatives. The gradient, or score, is

$$s_n = \nabla_{\beta} l_n = \sum_{i \in C} d_{in} \cdot [x_{in} - \sum_{k \in C} x_{kn} \cdot P_{Cn}(k)]$$

$$= \sum_{i \in C} [d_{in} - P_{Cn}(i)] \cdot x_{in} = \sum_{i \in C} [d_{in} - P_{cn}(i)] \cdot x_{iCn}$$

where
$$x_{Cn} = \sum_{i \in C} P_{Cn}(i) \cdot x_{in}$$
 and $x_{iCn} = x_{in} - x_{Cn}$

The score has the interpretation of requiring orthogonality in the sample between the explanatory variables x_{in} and the residuals d_{in} - $P_{Cn}(i)$. The hessian, or *information matrix*, is

$$H_{\rm n} = -\nabla_{\beta\beta} l_{\rm n} = \sum_{i \in C} P_{\rm Cn}(i) \cdot [x_{\rm in} - x_{\rm Cn}] \cdot [x_{\rm in} - x_{\rm Cn}]' = \sum_{i \in C} P_{\rm Cn}(i) \cdot x_{i{\rm Cn}} \cdot x_{i{\rm Cn}}',$$

The matrix H_n is positive semi-definite, and the expectation of H_n will be positive definite so long as the x_{iCn} are not linearly dependent. This assures that the log likelihood function is concave.

Consider the sample log likelihood $L_N = \sum_{n=1}^{N} l_n$. Any parameter vector that sets the sample

score to zero will also be a global maximum, and standard iterative maximization by a procedure like Newton-Raphson will converge to a global maximum.³ The Newton-Raphson iterative adjustment in parameters will be

$$\Delta\beta = \left(\sum_{n=1}^{N} H_{n}\right)^{-1} \sum_{n=1}^{N} s_{n} = \left(\sum_{n=1}^{N} \sum_{i \in C} P_{Cn}(i) x_{iCn} x_{iCn}'\right)^{-1} \sum_{n=1}^{N} \sum_{i \in C} x_{iCn} \cdot [d_{in} - P_{Cn}(i)],$$

where $x_{iCn} = x_{in} - x_{Cn} = x_{in} - \sum_{i \in C} P_{Cn}(i) \cdot x_{in}$. The adjustment $\Delta \beta$ can also be interpreted as the

estimates of the coefficients from a linear regression of $[d_{in} - P_{cn}(i)] \cdot P_{Cn}(i)^{-1/2}$ on the variables $P_{Cn}(i)^{1/2} \cdot x_{iCn}$. This has the same form as a Lagrange Multiplier test statistic, and one can write down a criterion for convergence that is identical to a LM test of whether the last iterate of the parameter vector is the true parameter vector. (One would want to accept the hypothesis and stop iterating only if there is very little probability of a type II error, accepting a false hypothesis. Therefore, the convergence criterion should use this LM statistic with a very *large* type I error, say 99.9%.)

One implication of the MNL model is that the ratio of the probabilities of two alternatives i and j depends only on x_i and x_j , and not on the presence or properties of other alternatives; i.e., $P_{Cn}(i)/P_{Cn}(j) = \exp((x_{in} - x_{jn})\beta)$. This is called the *Independence from Irrelevant Alternatives* (IIA)

³ A step size adjustment may speed convergence or avoid "overshooting" that could interfere with convergence.

property. This is a very restrictive property when x_{in} depends only on attributes of alternative i for each i. It implies patterns of cross-elasticities of substitution that are implausible for many applications. For example, a MNL model of the multinomial choice of school for graduate study in economics makes no allowance for the possibility that there may be unobserved factors shared by several schools (e.g., the Northern California location of Berkeley and Stanford), so that discrimination within this class (which we might call the "blue department" and the "red department") is likely to be sharper than it is between one of these departments and an East Coast department such as Princeton. The IIA property is a powerful restriction which if true can greatly simplify estimation and forecasting, and if false produces a misspecified model that can give misleading estimates and forecasts. The IIA property is not on its face particularly plausible, and what is remarkable about the MNL model is that it often performs well in forecasting situations even when IIA does not appear to be reasonable. However, it is important to understand the consequences of the IIA property of MNL, and to develop models for discrete response that can be used when IIA is clearly invalid.

7. ALTERNATIVES TO THE MNL MODEL FOR MULTINOMIAL RESPONSE

As in the derivation of the MNL model, associate with alternative i in a feasible set C a "payoff" $u_i = z_i \beta + \epsilon_i$, which in the case of consumer choice may be the indirect utility attached to alternative i and in the case of firm choice may be profit from alternative i. The z_i are observed explanatory variables, and the ϵ_i are unobserved disturbances. Observed choice is assumed to maximize payoff: $y_i = \mathbf{1}(u_i \ge u_j \text{ for } j \in C)$. One form of this model is a <u>random coefficients</u> formulation $u_i = z_i \alpha$, $\mathbf{E}\alpha = \beta$, $\epsilon_i = z_i (\alpha - \beta)$, implying $\text{cov}(\epsilon_i, \epsilon_j) = z_i \cdot \text{Cov}(\alpha) \cdot z_j'$. For $C = \{1, ..., J\}$, define u, z, ϵ , and y to be $J \times 1$ vectors with components $u_j, z_j, \epsilon_j, y_j$, respectively. Define a $(J-1) \times J$ matrix Δ_i by starting from the $J \times J$ identity matrix, deleting row i, and then replacing column i with the vector (-1, ..., -1). For example, letting $\mathbf{1}_{J-1}$ denote a $(J-1) \times 1$ vector of ones and $\mathbf{1}_{J-1}$ denote an identity matrix of dimension J-1, one has

$$\Delta_1 = [-1_{J-1} \ I_{J-1}].$$

Then alternative i is chosen if $\Delta_i u \leq 0$. The probability of this event is

$$P_i(z,\theta) = Pr(\Delta_i u \le 0 | z,\theta) = \int_{\Delta,u \le 0} f(u | z,\theta) du,$$

where $f(u|z,\theta)$ is the conditional density of u given z. The parameters θ include the slope parameters β and any additional parameters characterizing the distribution of the disturbances ϵ . The multivariate integral defining $P_i(z,\theta)$ can be calculated analytically in special cases, notably multinomial logit and its generalizations. However, for most densities the integral is analytically intractable, and for dimensions much larger than J=5 is also intractable to evaluate with adequate precision using standard numerical integration methods. Then, the four practical methods of working with random utility models for complex applications are (1) use of nested multinomial logit and related specializations of Generalized Extreme Value (GEV) models, (2) use of multinomial probit with special factor-analytic structure to provide feasible numerical integration; (3) use of

multinomial probit with simulation estimators that handle high dimensions; and (4) use of mixed (random coefficients) multinomial logit, with simulation procedures for the coefficients.

GEV Models

Assume that the indirect utility of i can be written $u_i = v_i + \varepsilon_i$ with ε_i a disturbance and v_i the systematic part of utility, depending on observed variables and unknown parameters. For example, one might have $v_i = \alpha(y-t_i) + \gamma x_i$, where y is income, t_i is the cost of alternative i (including costs of time), and ε_i is a part that varies randomly across consumers. The terms α , γ are parameters.

The E's have a joint CDF of generalized extreme value (GEV) form if

$$F(\varepsilon_1,...,\varepsilon_J) = \exp(-H(e^{-\varepsilon_1},...,e^{-\varepsilon_J})),$$

where (i) $H(w_1,...,w_J)$ is a non-negative linear homogeneous function of $w \ge 0$, satisfying (ii) if any argument goes to $+\infty$, then H goes to $+\infty$; and (iii) the mixed partial derivatives of H exist, are continuous, and alternate in sign, with non-negative odd mixed derivatives. A function H with properties (i) - (iii) will be termed a *GEV generating function*.

Theorem 1. Suppose H(w) for $w = (w_1,...,w_J)$ is a GEV generating function. Then, $F(\varepsilon)$ is a CDF with Extreme Value Type I univariate marginals. Further the random utility model $u_i = v_i + \varepsilon_i$ with ε distributed $F(\varepsilon)$ satisfies

E max_i
$$u_i = log H(e^{v_1}, ..., e^{v_J}) + E_i$$

where E = 0.5772156649 is Euler's constant, and the choice probabilities satisfy

$$P_{i} = e^{v_{i}} \cdot H_{i} (e^{v_{1}}, ..., e^{v_{J}})/H(e^{v_{1}}, ..., e^{v_{J}}).$$

The linear function $H = \sum_{i=1}^{J} w_i$ is a GEV generating function which yields the multinomial

logit (MNL) model. The following result can be used to build up complex choice models. In this theorem, the sets A and B are not required to be mutually exclusive.

Theorem 2. If sets A,B satisfy $A \cup B = \{1,...,J\}$, $H^A(w_A)$ and $H^B(w_B)$ are GEV generating functions in w_A and w_B , respectively, and if $s \ge 1$, then $H(w) = H^A(w_A^s)^{1/s} + H^B(w_B)$ is a GEV generating function in $(w_1,...,w_J)$.

One can use this theorem to show that a three-level nested MNL model is generated by a function H of the form

$$H = \sum_{m=1}^{M} \left[\sum_{k=1}^{K} \left[\sum_{i \in A_{mk}}^{\cdot} w_{i}^{s_{m} s_{k}'} \right] \frac{1}{s_{m}} \right]^{\frac{1}{s_{k}'}}$$

where the A_{mk} partition $\{1,...,J\}$ and $s'_k,s_m \ge 1$. This form corresponds to a tree: m indexes major branches, k indexes limbs from each branch, and i indexes the final twigs. The larger s'_k or s_m , the

more substitutable the alternatives in A_{mk} . If $s'_k = s_m = 1$, this model reduces to the MNL model. The GEV model is most efficiently estimated by MLE, but a convenient (and numerically relatively stable) method of getting preliminary estimates is to proceed sequentially, starting at the innermost nests. At each level of nesting, choice can be represented by a MNL model, which will however depend on parameters estimated from deeper levels of nesting. Details of this estimation procedure are given in McFadden (1984).

One interesting feature of GEV models is that they provide a convenient computational formula for the exact consumers' surplus associated with a policy that changes the attributes of alternatives. Let $v_i' = \alpha(y-t_i) + \gamma \ x_i'$ and $v_i'' = \alpha(y-t_i) + \gamma \ x_i''$, where x_i' is the vector of original attributes and x_i'' is the vector of improved attributes. Then, the willingness-to-pay for the change from x' to x'' is

WTP =
$$\frac{1}{\alpha} \cdot \left\{ \log H(e^{v''_1}, ..., e^{v''_J}) - \log H(e^{v'_1}, ..., e^{v'_J}) \right\}$$
.

This is the "log sum" formula first developed by Ben Akiva (1972), McFadden (1973), and Domencich and McFadden (1975) for the multinomial logit model, and by McFadden (1978, 1981) for the nested logit model. This formula is valid *only* when the indirect utility function is linear in income.

The MNP Model

A density that is relatively natural for capturing unobserved effects, and the patterns of correlation of these effects across alternatives, is the multivariate normal distribution with a flexible covariance matrix. This is termed the multinomial probit model. If $\varepsilon = z\xi$, where ξ is interpreted as a random variation in "taste" weights across observations with $\xi \sim N(0,\Omega)$, then the transformed variable $w = \Delta_i u$ is multivariate normal of dimension J-1 with mean $\Delta_i z \beta$ and covariance $\Delta_i z \Omega z' \Delta_i'$. Unless $J \leq 5$ or dimensionality can be reduced because ξ has a factorial covariance structure, the resulting MNP response probabilities are impractical to calculate by numerical integration. The method of simulated moments was initially developed to handle this model; see McFadden (1989).

For dynamic applications (e.g., multiperiod binomial probit with autocorrelation), and other applications with large dimension, alternatives to simulation of the MNP model with a unrestricted covariance matrix may perform better. McFadden (1984, 1989) suggests a "factor analytic" MNP with a components of variance structure, starting from

$$u_i = z_i \beta + \sum_{k=1}^{K} \lambda_{ik} \xi_k + \sigma_i \nu_i$$
,

where $\xi_1,...,\xi_K,v_1,...,v_J$ are independent standard normal, with the ξ_k interpreted as levels of unobserved factors and the λ_{ik} as the loading of factor k on alternative i. The λ 's are identified by normalizations and exclusion restrictions. The choice probabilities for this specification are

$$P_{i}(z,\theta) = \int_{\nu_{i}=-\infty}^{+\infty} \int_{\xi=-\infty}^{+\infty} \phi(\nu_{i}) \cdot \prod_{k=1}^{K} \phi(\xi_{k})$$

$$\times \prod_{j \neq i} \Phi \left(\frac{(z_j - z_i)\beta + \sum_k [\Lambda_{jk} - \Lambda_{ik}] \cdot \xi_k + \sigma_i v_i}{\sigma_j} \right) \cdot dv_i d\xi_1 \cdots d\xi_K$$

Numerical integration (when K+1 < 5) or simulation methods can be used to approximate this function and its derivatives for purposes of approximate maximum likelihood estimation. If simulation is used, two important rules should be followed: First, the Monte Carlo draws used for simulation should be made once and then frozen over the course of iterative search for parameters. This avoids "chatter" that can destroy the statistical properties of simulation-based estimators. Second, the number of simulation draws per observation should rise faster than the square root of sample size. This will assure that the simulation is asymptotically negligible, and cannot interfere with the CAN properties of MLE.

Mixed MNL (MMNL)

Mixed MNL is a generalization of standard MNL that shares many of the advantages of MNP, allowing a broad range of substitution patterns. Train and McFadden (1999) show that any regular random utility model can be approximated as closely as one wants by a MMNL model. Assume $u_i = z_i \alpha + \epsilon_i$, with the ϵ_i independently identically Extreme Value I distributed, and α random with density $f(\alpha;\theta)$, where θ is a vector of parameters. Conditioned on α ,

$$L_{i}(z|\alpha) = e^{z_{i}\alpha} / \sum_{i \in C} e^{z_{j}\alpha}$$
.

Unconditioning on α,

$$P_i(z|\theta) = \int_{\alpha} L_i(z|\alpha) f(\alpha;\theta) d\alpha$$
.

This model can be estimated by sampling randomly from $f(\alpha;\theta)$, approximating $P_i(z|\theta)$ by an average in this Monte Carlo sample, and varying θ to maximize the likelihood of the observations. Care must be taken to avoid chatter in the draws when θ varies. The MMNL model has proved computationally practical and flexible in applications. It can approximate MNP models well, and provides one convenient route to specification of models with flexibility comparable to that provided by MNP.

8. TESTS FOR THE IIA PROPERTY OF MNL

Alternatives to the MNL model may be derived from random utility models in which subsets of alternatives have disturbances ϵ_{in} that are correlated, perhaps because of common unobserved attributes. Common components of disturbances cancel out of the determination of choice within such a subset. As a result, discrimination of differences in observed attributes is sharper in a subset than overall; there is less random noise to blur discrimination. Tests for the presence of sharper discrimination in subsets is then a test of the IIA property of the MNL model.

For any discrete response model, including but not limited to MNL, let s_n denote the score of an observation, and H_n the negative of the hessian for an observation. A Taylor's expansion of the sample score about the maximum likelihood estimator establishes that in large samples

$$b - \beta_0 = (\sum_{n=1}^{N} H_n)^{-1} (\sum_{n=1}^{N} s_n) + O(N^{-1/2}),$$

and the covariance matrix of $b - \beta_0$ is approximately $\Omega = (\sum_{n=1}^{N} H_n)^{-1}$, where all expressions

are evaluated at β_o . In sufficiently large samples, b is approximately normally distributed with mean β_o and covariance matrix Ω , and the quadratic form

$$(b - \beta_0)'\Omega_C^{-1}(b - \beta_0) = (\sum_{n=1}^N s_n)'(\sum_{n=1}^N H_n)^{-1}(\sum_{n=1}^N s_n)$$

is approximately chi-squared distributed with degrees of freedom equal to the dimension of β_o . This is a Wald test statistic for the null hypothesis that $\beta = \beta_o$. It can also be applied to a subvector of β , with the commensurate submatrix of Ω_C^{-1} in the center of the quadratic form, to test the null hypothesis that this subvector takes on specified values.

We describe a series of hypothesis testing procedures that can be interpreted as tests of the IIA property of MNL. We will show a connection between these statistics and conventional test statistics for omitted variables.

Hausman-McFadden IIA Test:4

Estimate the MNL model twice, once on a full set of alternatives C, and second on a specified subset of alternatives A and the subsample with choices from this subset. If IIA holds, the two estimates should not be statistically different. If IIA fails, then there may be sharper discrimination within the subset A, so that the estimates from the second setup will be larger in magnitude than the estimates from the full set of alternatives. Let β_A denote the estimates obtained from the second setup, and Ω_A denote their estimated covariance matrix. Let β_C denote the estimated covariance matrix. (Some parameters that can be estimated from the full choice set may not be identified in the second setup, in which case β_C refers to estimates of the subvector of parameters that are identified in both setups.) Consider the quadratic form

$$(\beta_C$$
 - $\beta_A)'(\Omega_A$ - $\Omega_C)^{\text{--}1}(\beta_C$ - $\beta_A)$.

This has a chi-square distribution when IIA is true. In calculating this test, one must be careful to restrict the comparison of parameters, dropping components as necessary, to get $\Omega_{\rm A}$ - $\Omega_{\rm C}$ non-singular. When this is done, the degrees of freedom of the chi-square test equals the rank of $\Omega_{\rm A}$ - $\Omega_{\rm C}$. The simple form of the covariance matrix for the parameter difference arises because $\beta_{\rm C}$ is the efficient estimator for the problem.

⁴Hausman-McFadden, *Econometrica*, 1984.

McFadden omitted variables test.⁵

Estimate the basic MNL model, using all the observations; let $P_{\rm in} = P_{\rm Cn}(i)$ denote the fitted model. Suppose A is a specified subset of alternatives. Create new variables in one of the following three forms:

a. If x_{in} are the variables in the basic logit model, define new variables

$$z_{in} = \begin{cases} x_{in} - (\sum_{j \in A} P_{jn} x_{jn}) / (\sum_{j \in A} P_{jn}) & if i \in A \\ 0 & if i \notin A \end{cases},$$

The variables z_{in} can be written in abbreviated form as $z_{in} = \delta_{iA}(x_{in} - x_{An})$, where $\delta_{iA} = 1$ iff

 $i \in A$ and $x_{An} = \sum_{j \in A} P_{jn|A} x_j$ and $P_{jn|A}$ is the conditional probability of choice of j given

choice from A, calculated from the base model.

b. If $V_{in} = x_{in}\beta$ is the representative utility from the basic model, calculated at basic model estimated parameters, define the new variable

$$\mathbf{z}_{\text{in}} = \begin{cases} V_{in} - (\sum_{j \in A} P_{jn} V_{jn}) / (\sum_{j \in A} P_{jn}) & \text{if } i \in A \\ 0 & \text{if } i \notin A \end{cases},$$

or more compactly, $z_{in} = \delta_{iA}(V_{in} - V_{An})$.

c. Define the new variable

$$\mathbf{z}_{\text{in}} = \begin{cases} \log(P_{in|A}) - \sum_{k \in A} P_{kn|A} \log(P_{kn|A}) & \text{if } i \in A \\ 0 & \text{if } i \notin A \end{cases},$$

where $P_{\text{in}|A}$ is calculated using the basic model estimates.

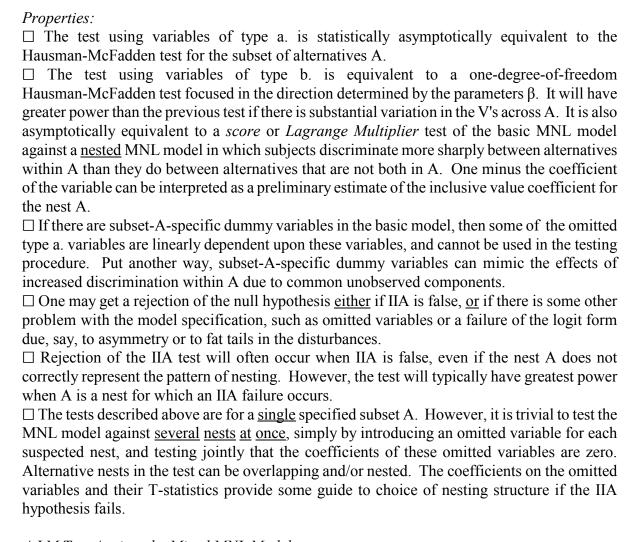
- \square The constructions b. and c. are the <u>same</u>. The denominators of the probabilities in the expression $-log(P_{in|A})$ that appears in the type c. variable drop out, leaving the terms in the construction b.
- \square Estimate an expanded MNL model that contains the basic model variables plus the new variables z_{in} . Then test whether these added variables are significant. If there is a single added

⁵D. McFadden, "Regression based specification tests for the multinomial logit model" *Journal of Econometrics*, 1987.

variable, as in the construction b., then the T-statistic for this added variable is the appropriate test statistic. More generally, one can form a likelihood ratio statistic

$$LR = 2 \left[\begin{pmatrix} Log \ Likelihood \\ with \ z's \end{pmatrix} - \begin{pmatrix} Log \ Likelihood \\ without \ z's \end{pmatrix} \right]$$

If IIA holds, this likelihood ratio statistic has a chi-square distribution with degrees of freedom equal to the number of added z variables (after eliminating any that are linearly dependent).



A LM Test Against the Mixed MNL Model

The mixed MNL family is very flexible and can approximate any well-behaved discrete response data generation process that is consistent with utility maximization. However, because the MMNL model requires the use of simulation methods for estimation, it is very useful to have a specification test that can indicate whether mixing is needed. The next result describes

a Lagrange Multiplier test for this purpose. This test has the pivotal property that its asymptotic distribution under the null hypothesis that the correct specification is MNL does not depend on the parameterization of the mixing distribution under the alternative.

Theorem 2. Consider choice from a set $\mathbb{C} = \{1,...,J\}$. Let x_i be a $1 \times K$ vector of attributes of alternative i. From a random sample n = 1,...,N, estimate the parameter α in the simple MNL

 $model\ L_{\mathbf{C}}(\mathbf{i};\mathbf{x},\alpha) = e^{x_i\alpha}/\sum_{j\in\mathbf{C}}e^{x_j\alpha}$, using maximum likelihood; construct artificial variables

$$z_{ti} = \frac{1}{2}(x_{ti} - x_{tC})^2$$
 with $x_{tC} = \sum_{j \in C} x_{tj} \cdot L_C(j; \mathbf{x}, \hat{\alpha})$

for selected components t of x_i , and use a Wald or Likelihood Ratio test for the hypothesis that the artificial variables z_{ti} should be omitted from the MNL model. This test is asymptotically equivalent to a Lagrange multiplier test of the hypothesis of no mixing against the alternative of a MMNL model $P_C(i|x,\theta) = \int L_C(i;x,\alpha) \cdot G(d\alpha;\theta)$ with mixing in the selected components t of α . The degrees of freedom equals the number of artificial variables z_{ti} that are linearly independent of x.

To examine the operating characteristics of the test, consider two simple Monte Carlo experiments for choice among three alternatives, with random utility functions $u_i = \alpha_1 x_{1i} + \alpha_2 x_{2i} + \epsilon_i$. The disturbances ϵ_i were i.i.d. Extreme Value Type I. In the first experiment, the covariate were distributed as described below:

Variable	Alternative 1	Alternative 2	Alternative 3
\mathbf{x}_1	$\pm \frac{1}{2}$ w.p. $\frac{1}{2}$	0	0
\mathbf{x}_2	$\pm \frac{1}{2}$ w.p. $\frac{1}{2}$	$\pm \frac{1}{2}$ w.p. $\frac{1}{2}$	0

The parameter $\alpha_2 = 1$ under both the null and the alternative. The parameter $\alpha_1 = 0.5$ under the null hypothesis, and under the alternative $\alpha_1 = 0.5 \pm 1$ w.p. 1/2. We carried out 1000 repetitions of the test procedure for a sample of size N = 1000 and choices generated alternately under the null hypothesis and under the alternative just described, using likelihood ratio tests for the omitted variable z_{1i} . The results are given below:

Nominal Significance Level	Actual Significance Level	Power Against the Alternative
10%	8.2%	15.6%
5%	5.0%	8.2%

The nominal and actual significance levels of the test agree well. The power of the test is low, and an examination of the estimated coefficients reveals that the degree of heterogeneity in tastes present in this experiment gives estimated coefficients close to their expected values. Put another way, this pattern of heterogeneity is difficult to distinguish from added extreme value noise.

In the second experiment, the covariates are distributed as shown below:

Variable	Alternative 1	Alternative 2	Alternative 3
\mathbf{X}_1	$\pm \frac{1}{2}$ w.p. $\frac{1}{2}$	$\pm \frac{1}{2}$ w.p. $\frac{1}{2}$	0
\mathbf{X}_2	$\pm \frac{1}{2}$ W.p. $\frac{1}{2}$	$\pm \frac{1}{2}$ w.p. $\frac{1}{2}$	0

The utility function is again $u_i = \alpha_1 x_{1i} + \alpha_2 x_{2i} + \epsilon_i$. Under the null hypothesis, $\alpha_1 = \alpha_2 = 1$, while under the alternative $(\alpha_1, \alpha_2) = (2,0)$ w.p. ½ and (0,2) w.p. ½. Again, 1000 repetitions of the tests are made for N = 1000 under the null and the alternative; the results are given below:

Nominal Significance Level	Actual Significance Level	Power Against the Alternative
10%	9.7%	52.4%
5%	3.9%	39.8%

In this case where mixing is across utility functions of different variables, the test is moderately powerful. It remains the case in this example that the estimated coefficients in the MNL model without mixing are close to their expected values.