



$$y_i = \begin{cases} 1, & \text{if } y_i^* > 0, \\ 0, & \text{if } y_i^* \leq 0. \end{cases} \quad (6.2.2)$$

The expected value of  $y_i$  is then the probability that the event will occur,

$$\begin{aligned} E(y_i | \mathbf{x}_i) &= 1 \cdot Pr(v_i > -\beta' \mathbf{x}_i) + 0 \cdot Pr(v_i \leq -\beta' \mathbf{x}_i) \\ &= Pr(v_i > -\beta' \mathbf{x}_i) \\ &= Pr(y_i = 1 | \mathbf{x}_i) \end{aligned} \quad (6.2.3)$$

When the probability law of generating  $v_i$  follows a two-point distribution  $(1 - \beta' \mathbf{x}_i)$  and  $(-\beta' \mathbf{x}_i)$ , with probabilities  $\beta' \mathbf{x}_i$  and  $(1 - \beta' \mathbf{x}_i)$ , respectively, we have the linear-probability model

$$y_i = \beta' \mathbf{x}_i + v_i \quad (6.2.4)$$

with  $E v_i = \beta' \mathbf{x}_i (1 - \beta' \mathbf{x}_i) + (1 - \beta' \mathbf{x}_i)(-\beta' \mathbf{x}_i) = 0$ . When the probability density function of  $v_i$  is a standard normal density function,  $\frac{1}{\sqrt{2\pi}} \exp(-\frac{v^2}{2}) = \phi(v)$ , we have the Probit model,

$$\begin{aligned} Pr(y_i = 1 | \mathbf{x}_i) &= \int_{-\beta' \mathbf{x}_i}^{\infty} \phi(v_i) dv_i \\ &= \int_{-\infty}^{\beta' \mathbf{x}_i} \phi(v_i) dv_i = \Phi(\beta' \mathbf{x}_i). \end{aligned} \quad (6.2.5)$$

When the probability density function is a standard logistic,

$$\frac{\exp(v_i)}{(1 + \exp(v_i))^2} = [(1 + \exp(v_i))(1 + \exp(-v_i))]^{-1}$$

we have the logit model

$$Pr(y_i = 1 | \mathbf{x}_i) = \int_{-\beta' \mathbf{x}_i}^{\infty} \frac{\exp(v_i)}{(1 + \exp(v_i))^2} dv_i = \frac{\exp(\beta' \mathbf{x}_i)}{1 + \exp(\beta' \mathbf{x}_i)}. \quad (6.2.6)$$

Let  $F(\beta' \mathbf{x}_i) = E(y_i | \mathbf{x}_i)$ . The three commonly used parametric models for the binary choice may be summarized with a single index  $w$  as:

*Linear-Probability Model*

$$F(w) = w. \quad (6.2.7)$$

*Probit model*

$$F(w) = \int_{-\infty}^w \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du = \Phi(w). \quad (6.2.8)$$

*Logit model*

$$F(w) = \frac{e^w}{1 + e^w}. \quad (6.2.9)$$

The linear-probability model is a special case of the linear regression model with heteroscedastic variance,  $\beta' \mathbf{x}_i (1 - \beta' \mathbf{x}_i)$ . It can be estimated by least squares or weighted least squares (Goldberger 1964). But it has an obvious defect in that  $\beta' \mathbf{x}_i$  is not constrained to lie between 0 and 1 as a probability should, whereas the Probit and logit models do. The probability functions used for the Probit and logit models are the standard normal distribution and the logistic distribution, respectively. We use cumulative standardized distribution because in the dichotomy case there is no information about the spread of

a continuous random variable. For instance, for the Probit model, the probability that an event occurs depends only on  $(\frac{1}{\sigma})\beta'x$ , where  $\sigma$  denotes the standard deviation of a normal density. There is no way to identify the variance of a normal density. The logit probability density function is symmetric around 0 and has a variance of  $\pi^2/3$ . Because they are distribution functions, the Probit and logit models are bounded between 0 and 1.

The cumulative normal distribution and the logistic distribution are very close to each other, except that the logistic distribution has slightly heavier tails (Cox 1970). Moreover, the cumulative normal distribution  $\Phi$  is reasonably well approximated by a linear function for the range of probabilities between 0.3 and 0.7. Amemiya (1981) has suggested an approximate conversion rule for the coefficients of these three models. Let the coefficients for the linear-probability, Probit, and logit models be denoted as  $\hat{\beta}_{LP}$ ,  $\hat{\beta}_{\Phi}$ , and  $\hat{\beta}_L$ , respectively. Then

$$\hat{\beta}_L \simeq 1.6 \hat{\beta}_{\Phi},$$

$$\hat{\beta}_{LP} \simeq 0.4 \hat{\beta}_{\Phi} \text{ except for the constant term,} \quad (6.2.10)$$

and

$$\hat{\beta}_{LP} \simeq 0.4 \hat{\beta}_{\Phi} + 0.5 \text{ for the constant term.}$$

For a random sample of  $N$  individuals, the likelihood function for these three models can be written in general form as

$$L = \prod_{i=1}^N F(\beta'x_i)^{y_i} [1 - F(\beta'x_i)]^{1-y_i}. \quad (6.2.11)$$

Differentiating the logarithm of the likelihood function yields the vector of first derivatives and the matrix of second-order derivatives as

$$\frac{\partial \log L}{\partial \beta} = \sum_{i=1}^N \frac{y_i - F(\beta'x_i)}{F(\beta'x_i)[1 - F(\beta'x_i)]} F'(\beta'x_i)x_i \quad (6.2.12)$$

$$\begin{aligned} \frac{\partial \log L}{\partial \beta \partial \beta'} = \sum_{i=1}^N \left\{ - \left[ \frac{y_i}{F^2(\beta'x_i)} + \frac{1 - y_i}{[1 - F(\beta'x_i)]^2} \right] [F'(\beta'x_i)]^2 \right. \\ \left. + \left[ \frac{y_i - F(\beta'x_i)}{F(\beta'x_i)[1 - F(\beta'x_i)]} \right] F''(\beta'x_i) \right\} x_i x_i' \end{aligned} \quad (6.2.13)$$

where  $F'(\beta'x_i)$  and  $F''(\beta'x_i)$  denote the first and second derivatives of  $F(\beta'x_i)$  with respect to  $\beta'x_i$ . If the likelihood function (6.2.11) is concave, as in the models discussed here (e.g., Amemiya 1985, p. 273), a Newton–Raphson iterative method,

$$\hat{\beta}^{(j)} = \hat{\beta}^{(j-1)} - \left( \frac{\partial^2 \log L}{\partial \beta \partial \beta'} \right)_{\beta=\hat{\beta}^{(j-1)}}^{-1} \left( \frac{\partial \log L}{\partial \beta} \right)_{\beta=\hat{\beta}^{(j-1)}}, \quad (6.2.14)$$

or a method of scoring,

$$\hat{\beta}^{(j)} = \hat{\beta}^{(j-1)} - \left[ E \frac{\partial^2 \log L}{\partial \beta \partial \beta'} \right]_{\beta=\hat{\beta}^{(j-1)}}^{-1} \left( \frac{\partial \log L}{\partial \beta} \right)_{\beta=\hat{\beta}^{(j-1)}}, \quad (6.2.15)$$

can be used to find the maximum-likelihood estimator of  $\beta$ , where  $\hat{\beta}^{(j)}$  denotes the  $j$ th iterative solution.



If the outcomes are unordered, for instance,

$$y_i = \begin{cases} 1, & \text{if mode of transport is car,} \\ 2, & \text{if mode of transport is bus,} \\ 3, & \text{if mode of transport is train,} \end{cases}$$

then we will have to use a multivariate probability distribution to characterize the outcomes. One way to postulate unordered outcomes is to assume that the  $j$ th alternative is chosen because it yields higher utility than the utility of other alternatives. Let the  $i$ th individual's utility of choosing  $j$ th alternative be

$$y_{ij}^* = \mathbf{x}_i' \boldsymbol{\beta}_j + v_{ij}, j = 0, 1, \dots, m_i. \quad (6.2.19)$$

Then

$$\begin{aligned} \text{Prob}(y_i = j \mid \mathbf{x}_i) &= \text{Prob}(y_{ij}^* > y_{i\ell}^*, \quad \forall \ell \neq j \mid \mathbf{x}_i) \\ &= F_{ij}. \end{aligned} \quad (6.2.20)$$

The probability  $F_{ij}$  is derived from the joint distribution of  $(v_{i0}, \dots, v_{im})$ . If  $(v_{i0}, \dots, v_{im})$  follows a multivariate normal distribution, then (6.2.20) yields a multivariate Probit. If the errors  $v_{ij}$  are independently identically distributed with type I extreme value distribution, (6.2.20) yields a conditional logit model (McFadden 1974). However, contrary to the univariate case, the similar predictions in terms of the Probit and logit specifications no longer holds. In general, they will lead to different inferences. The advantage of the multivariate Probit model is that it allows the choice among alternatives to have an arbitrary correlation. The disadvantage is that the evaluation of  $\text{Prob}(y_i = j)$  involves multiple integrations which can be computationally infeasible. The advantage of the conditional logit model is that the evaluation of  $\text{Prob}(y_i = j)$  does not involve multiple integration. The disadvantage is that the relative odds between two alternatives are independent of the presence or absence of the other alternatives, the so-called *independence of irrelevant alternatives*. If the errors among alternatives are not independently distributed, this can lead to grossly false predictions of the outcomes. For discussion of model specification tests, see Hausman and McFadden (1984), Hsiao (1992b), Lee (1982, 1987), and Small and Hsiao (1985).

Because in many cases, a multi-response model can be transformed into a dichotomous model characterized by the  $\sum_{i=1}^N (m_i + 1)$  binary variables as in (6.2.16),<sup>2</sup> for ease of exposition, we shall concentrate only on the dichotomous model.<sup>3</sup>

When there is no information about the probability laws of generating  $v_i$ , a semiparametric approach can be used to estimate  $\boldsymbol{\beta}$  subject to the certain normalization rule (e.g., Klein and Spady 1993; Manski 1985; Powell, Stock, and Stoker 1989). (Section 6.4 will discuss further the semiparametric approach for panel data single index models.) However, whether an investigator takes a parametric or semiparametric approach, the cross-sectional model assumes that the error term  $v_i$  in the latent response function (6.2.1) is independently identically distributed and is independent of  $\mathbf{x}_i$ . In other words, conditional on  $\mathbf{x}_i$ , everyone has the same probability that an event will occur. It does not allow the possibility that the average behavior given  $\mathbf{x}$  can be different from individual probabilities; that is, it does not

<sup>2</sup> The variable  $y_{i0}$  is sometimes omitted from the specification because it is determined by  $y_{i0} = 1 - \sum_{j=1}^m y_{ij}$ . For instance, a dichotomous model is often simply characterized by a single binary variable  $y_i, i = 1, \dots, N$ .

<sup>3</sup> It should be noted that in generalizing the results of the binary case to the multiresponse case, we should allow for the fact that although  $y_{ij}$  and  $y_{i'j}$  are independent for  $i \neq i'$ ,  $y_{ij}$  and  $y_{ij'}$  are not, because  $\text{Cov}(y_{ij}, y_{ij'}) = -F_{ij}F_{ij'}$ .

allow  $Pr(y_i = 1 \mid \mathbf{x}) \neq Pr(y_j = 1 \mid \mathbf{x})$ . The availability of panel data provides the possibility of distinguishing average behavior from individual behavior by decomposing the error term,  $v_{it}$ , into

$$v_{it} = \alpha_i + \lambda_t + u_{it}, \quad (6.2.21)$$

where  $\alpha_i$  and  $\lambda_t$  denote the effects of omitted individual-specific and time-specific variables, respectively. Then  $\text{Prob}(y_i = 1 \mid \mathbf{x}, \alpha_i) \neq \text{Prob}(y_j = 1 \mid \mathbf{x}, \alpha_j)$  if  $\alpha_i \neq \alpha_j$ . In this chapter, we shall demonstrate the misspecifications that can arise because of failure to control for unobserved characteristics of the individuals in panel data and discuss possible remedies.

### 6.3 PANEL PARAMETRIC APPROACH TO STATIC MODELS WITH HETEROGENEITY

Statistical models developed for analyzing cross-sectional data essentially ignore individual differences and treat the sum of the individual-specific effect and the time-varying omitted-variable effect as a pure chance event. However, as the example in Chapter 1 shows, a discovery of a group of married women having an average yearly labor participation rate of 50% could lead to diametrically opposite inferences. At one extreme, each woman in a homogeneous population could have a 50% chance of being in the labor force in any given year; whereas at the other extreme, 50% of women in a heterogeneous population might always work and 50% never work. Either explanation is consistent with the finding relying on given cross-sectional data. To discriminate among the many possible explanations, we need information on individual labor-force histories in different subintervals of a life cycle. Panel data, having information on inter-temporal dynamics of individual entities, provide the possibility of separating a model of individual behavior from a model of average behavior of a group of individuals.

Suppose there are sample observations  $(y_{it}, \mathbf{x}_{it})$  for  $i = 1, \dots, N$  and  $t = 1, \dots, T$ , where  $y_{it}$  is binary with  $y_{it} = 1$  if  $y_{it}^*$  given by (6.2.1) is greater than 0, and 0 otherwise. For simplicity, we shall assume that the heterogeneity across cross-sectional units is time invariant,<sup>4</sup> and these individual-specific effects are captured by decomposing the error term  $v_{it}$  in (6.2.1) as  $\alpha_i + u_{it}$ . When  $\alpha_i$  are treated as fixed,  $\text{Var}(v_{it} \mid \alpha_i) = \text{Var}(u_{it}) = \sigma_u^2$ . When  $\alpha_i$  are treated as random, we assume that  $E\alpha_i = E\alpha_i u_{it} = 0$ , and  $\text{Var}(v_{it}) = \sigma_u^2 + \sigma_\alpha^2$ . However, as discussed earlier, when the dependent variables are binary, the scale factor is not identifiable. Thus, for ease of exposition, we normalize the variance of  $u, \sigma_u^2$ , to be equal to 1 for the parametric specifications discussed in Section 6.2.

The existence of such unobserved permanent components allows individuals who are homogeneous in terms of their observed characteristics to be heterogeneous in response probabilities,  $F(\beta' \mathbf{x}_{it} + \alpha_i)$ . For example, heterogeneity will imply that the sequential-participation behavior of a woman,  $F(\beta' \mathbf{x} + \alpha_i)$ , within a group of women with observationally identical  $\mathbf{x}$  differ systematically from  $F(\beta' \mathbf{x})$  or the average behavior of the group,  $\int F(\beta' \mathbf{x} + \alpha) dH(\alpha \mid \mathbf{x})$ , where  $H(\alpha \mid \mathbf{x})$  gives the population probability (or empirical distribution) for  $\alpha$  conditional on  $\mathbf{x}$ .<sup>5</sup> In this section, we discuss statistical inference of the common parameters  $\beta$  based on a parametric specification of  $F(\cdot)$ .

<sup>4</sup> For a random-coefficient formulation of Probit models, see Hausman and Wise (1978).

<sup>5</sup> Note that, in general,  $\int F(\beta' \mathbf{x} + \alpha) dH(\alpha \mid \mathbf{x}) \neq F[\beta' \mathbf{x} + E(\alpha \mid \mathbf{x})]$ .



Solving (6.3.4), we have

$$\begin{aligned}\hat{\alpha}_i &= \infty \text{ if } y_{i1} + y_{i2} = 2, \\ \hat{\alpha}_i &= -\infty \text{ if } y_{i1} + y_{i2} = 0, \\ \hat{\alpha}_i &= -\frac{\beta}{2} \text{ if } y_{i1} + y_{i2} = 1.\end{aligned}\tag{6.3.5}$$

Inserting (6.3.5) into (6.3.3) and letting  $n_1$  denote the number of individuals with  $y_{i1} + y_{i2} = 1$  and letting  $n_2$  denote the number of individuals with  $y_{i1} + y_{i2} = 2$ , we have<sup>7</sup>

$$\sum_{i=1}^N \frac{e^{\beta+\alpha_i}}{1 + e^{\beta+\alpha_i}} = n_1 \frac{e^{\beta/2}}{1 + e^{\beta/2}} + n_2 = \sum_{i=1}^N y_{i2}.\tag{6.3.6}$$

Therefore,

$$\hat{\beta} = 2 \left\{ \log \left( \sum_{i=1}^N y_{i2} - n_2 \right) - \log \left( n_1 + n_2 - \sum_{i=1}^N y_{i2} \right) \right\}.\tag{6.3.7}$$

By a law of large numbers (Rao 1973, chapter 2),

$$\begin{aligned}\text{plim}_{N \rightarrow \infty} \frac{1}{N} \left( \sum_{i=1}^N y_{i2} - n_2 \right) \\ &= \frac{1}{N} \sum_{i=1}^N \text{Prob}(y_{i1} = 0, y_{i2} = 1 \mid \beta, \alpha_i) \\ &= \frac{1}{N} \sum_{i=1}^N \frac{e^{\beta+\alpha_i}}{(1 + e^{\alpha_i})(1 + e^{\beta+\alpha_i})},\end{aligned}\tag{6.3.8}$$

$$\begin{aligned}\text{plim}_{N \rightarrow \infty} \frac{1}{N} \left( n_1 + n_2 - \sum_{i=1}^N y_{i2} \right) \\ &= \frac{1}{N} \sum_{i=1}^N \text{Prob}(y_{i1} = 1, y_{i2} = 0 \mid \beta, \alpha_i) \\ &= \frac{1}{N} \sum_{i=1}^N \frac{e^{\alpha_i}}{(1 + e^{\alpha_i})(1 + e^{\beta+\alpha_i})}.\end{aligned}\tag{6.3.9}$$

Substituting  $\hat{\alpha}_i = \frac{\beta}{2}$  into (6.3.8) and (6.3.9) yields

$$\text{plim}_{N \rightarrow \infty} \hat{\beta} = 2\beta,\tag{6.3.10}$$

which is not consistent.

<sup>7</sup> The number of individuals with  $y_{i1} + y_{i2} = 0$  is  $N - n_1 + n_2$ .



## 6.3.1.2 Conditions for the Existence of a Consistent Estimator

Neyman and Scott (1948) have suggested a general principle to find a consistent estimator for the (structural) parameter  $\beta$  in the presence of the incidental parameters  $\alpha_i$ .<sup>8</sup> Suppose the dimension of  $\beta$  is  $K$ ; their idea is to find  $K$  functions

$$\Psi_{Nj}(\mathbf{y}_1, \dots, \mathbf{y}_N | \beta), \quad j = 1, \dots, K. \quad (6.3.11)$$

that are independent of the incidental parameters  $\alpha_i$  and have the property that when  $\beta$  are the true values,  $\Psi_{Nj}(\mathbf{y}_1, \dots, \mathbf{y}_N | \beta)$  converges to some known constant, say, zero, in probability as  $N$  tends to infinity. Then an estimator  $\hat{\beta}$  derived by solving  $\Psi_{Nj}(\mathbf{y}_1, \dots, \mathbf{y}_N | \hat{\beta}) = 0$  is consistent under suitable regularity conditions. For instance,  $\hat{\beta}^* = (1/2)\hat{\beta}$  for the foregoing example of a fixed-effect logit model is such an estimator.

In the case of a linear-probability model, either taking first difference over time or taking difference with respect to the individual mean eliminates the individual-specific effect. The least-squares regression of the differenced equations yields a consistent estimator for  $\beta$  when  $N$  tends to infinity.

But in the general nonlinear models, simple forms of  $\Psi(\cdot)$  are not always easy to find. For instance, in general, we do not know the probability limit of the MLE of a fixed-effects logit model. However, if a minimum sufficient statistic  $\tau_i$  for the incidental parameter  $\alpha_i$  exists and is not dependent on the structural parameter  $\beta$ , the conditional density,

$$f^*(\mathbf{y}_i | \beta, \tau_i) = \frac{f(\mathbf{y}_i | \beta, \alpha_i)}{g(\tau_i | \beta, \alpha_i)} \text{ for } g(\tau_i | \beta, \alpha_i) > 0, \quad (6.3.12)$$

no longer depends on  $\alpha_i$ .<sup>9</sup> Maximizing the conditional density of  $\mathbf{y}_1, \dots, \mathbf{y}_N$ , given  $\tau_1, \dots, \tau_N$ ,

$$\prod_{i=1}^N f^*(\mathbf{y}_i | \beta, \tau_i), \quad (6.3.13)$$

yields the first-order conditions  $\Psi_{Nj}(\mathbf{y}_1, \dots, \mathbf{y}_N | \hat{\beta}, \tau_1, \tau_2, \dots, \tau_N) = 0$ , for  $j = 1, \dots, K$ . E.B. Anderson (1970, 1973) has shown that solving these functions will give a consistent estimator of the common (structural) parameter  $\beta$  under mild regularity conditions.<sup>10</sup>

To illustrate the conditional maximum likelihood method, we use the logit model as an example. The joint probability of  $\mathbf{y}_i$  is

$$\text{Prob}(\mathbf{y}_i) = \frac{\exp\{\alpha_i \sum_{t=1}^T y_{it} + \beta' \sum_{t=1}^T \mathbf{x}_{it} y_{it}\}}{\prod_{t=1}^T [1 + \exp(\beta' \mathbf{x}_{it} + \alpha_i)]}. \quad (6.3.14)$$

<sup>8</sup> We call  $\beta$  the structural parameter because the value of  $\beta$  characterizes the structure of the complete sequence of random variables. It is the same for all  $i$  and  $t$ . We call  $\alpha_i$  an incidental parameter to emphasize that the value of  $\alpha_i$  changes when  $i$  changes.

<sup>9</sup> Suppose that the observed random variables  $\mathbf{y}$  have a certain joint distribution function that belongs to a specific family  $\mathcal{J}$  of distribution functions. The statistic  $S(\mathbf{y})$  (a function of the observed sample values  $\mathbf{y}$ ) is called a sufficient statistic if the conditional expectation of any other statistic  $H(\mathbf{y})$ , given  $S(\mathbf{y})$ , is independent of  $\mathcal{J}$ . A statistic  $S^*(\mathbf{y})$  is called a minimum sufficient statistic if it is a function of every sufficient statistic  $S(\mathbf{y})$  for  $\mathcal{J}$ . For additional discussion, see Zacks (1971, chapter 2).

<sup>10</sup> When  $u_{it}$  are independently normally distributed, the LSDV estimator of  $\beta$  for the linear static model is the conditional MLE (Cornwell and Schmidt 1984).

$\sum_{t=1}^T y_{it}$  is

$$\text{Prob} \left( \mathbf{y}_i \mid \sum_{t=1}^T y_{it} \right) = \frac{\exp \left[ \boldsymbol{\beta}' \sum_{t=1}^T \mathbf{x}_{it} y_{it} \right]}{\sum_{D_{ij} \in \bar{B}_i} \exp \{ \boldsymbol{\beta}' \sum_{t=1}^T \mathbf{x}_{it} d_{ijt} \}}, \quad (6.3.15)$$

where  $\bar{B}_i = \{D_{ij} = (d_{ij1}, \dots, d_{ijT}) \mid d_{ijt} = 0 \text{ or } 1, \text{ and } \sum_{t=1}^T d_{ijt} = \sum_{t=1}^T y_{it}, j = 1, 2, \dots, \frac{T!}{s!(T-s)!}\}$ , is the set of all possible distinct sequence  $(d_{ij1}, d_{ij2}, \dots, d_{ijT})$  satisfying  $\sum_{t=1}^T d_{ijt} = \sum_{t=1}^T y_{it} = s$ . There are  $T + 1$  distinct alternative sets corresponding to  $\sum_{t=1}^T y_{it} = 0, 1, \dots, T$ . Groups for which  $\sum_{t=1}^T y_{it} = 0$  or  $T$  contribute zero to the likelihood function, because the corresponding probability in this case is equal to 1 (with  $\alpha_i = -\infty$  or  $\infty$ ). So only  $T - 1$  alternative sets are relevant. The alternative sets for groups with  $\sum_{t=1}^T y_{it} = s$  have  $\binom{T}{s}$  elements, corresponding to the distinct sequences of  $T$  trials with  $s$  success.

Equation (6.3.15) is in a conditional logit form (McFadden 1974), with the alternative sets  $(\bar{B}_i)$  varying across observations  $i$ . It does not depend on the incidental parameters,  $\alpha_i$ . Therefore, the conditional MLE of  $\beta$  is consistent under mild conditions. For example, with  $T = 2$ , the only case of interest is  $y_{i1} + y_{i2} = 1$ . The two possibilities are  $\omega_i = 1$ , if  $(y_{i1}, y_{i2}) = (0, 1)$ , and  $\omega_i = 0$ , if  $(y_{i1}, y_{i2}) = (1, 0)$ .

The conditional probability of  $\omega_i = 1$  given  $y_{i1} + y_{i2} = 1$  is

$$\begin{aligned} \text{Prob}(\omega_i = 1 \mid y_{i1} + y_{i2} = 1) &= \frac{\text{Prob}(\omega_i = 1)}{\text{Prob}(\omega_i = 1) + \text{Prob}(\omega_i = 0)} \\ &= \frac{\exp[\beta'(\mathbf{x}_{i2} - \mathbf{x}_{i1})]}{1 + \exp[\beta'(\mathbf{x}_{i2} - \mathbf{x}_{i1})]} \\ &= F[\beta'(\mathbf{x}_{i2} - \mathbf{x}_{i1})]. \end{aligned} \quad (6.3.16)$$

Equation (6.3.16) is in the form of a binary logit function in which the two outcomes are (0,1) and (1,0), with explanatory variables  $(\mathbf{x}_{i2} - \mathbf{x}_{i1})$ . The conditional log-likelihood function is

$$\begin{aligned} \log L^* &= \sum_{i \in \tilde{B}_1} \{\omega_i \log F[\boldsymbol{\beta}'(\mathbf{x}_{i2} - \mathbf{x}_{i1})] \\ &\quad + (1 - \omega_i) \log (1 - F[\boldsymbol{\beta}'(\mathbf{x}_{i2} - \mathbf{x}_{i1}))]\}, \end{aligned} \quad (6.3.17)$$

where  $\tilde{B}_1 = \{i \mid y_{i1} + y_{i2} = 1\}$ .

Although  $\tilde{B}_1$  is a random set of indices, Chamberlain (1980) has shown that the inverse of the information matrix based on the conditional-likelihood function provides an asymptotic covariance matrix for the conditional MLE of  $\boldsymbol{\beta}$  when  $N$  tends to infinity. This can be made more explicit by defining  $d_i = 1$ , if  $y_{i1} + y_{i2} = 1$ , and  $d_i = 0$ , otherwise, for the foregoing case in which  $T = 2$ . Then we have

$$J_{\tilde{B}_1} = \frac{\partial^2 \log L^*}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = - \sum_{i=1}^N d_i F[\boldsymbol{\beta}'(\mathbf{x}_{i2} - \mathbf{x}_{i1})] \{1 - F[\boldsymbol{\beta}'(\mathbf{x}_{i2} - \mathbf{x}_{i1})]\} (\mathbf{x}_{i2} - \mathbf{x}_{i1}) \cdot (\mathbf{x}_{i2} - \mathbf{x}_{i1})'. \quad (6.3.18)$$

The information matrix is

$$J = E(J_{\tilde{B}_1}) = - \sum_{i=1}^N P_i F[\boldsymbol{\beta}'(\mathbf{x}_{i2} - \mathbf{x}_{i1})] \{1 - F[\boldsymbol{\beta}'(\mathbf{x}_{i2} - \mathbf{x}_{i1})]\} (\mathbf{x}_{i2} - \mathbf{x}_{i1}) \cdot (\mathbf{x}_{i2} - \mathbf{x}_{i1})', \quad (6.3.19)$$

where  $P_i = E(d_i \mid \alpha_i) = F(\boldsymbol{\beta}'\mathbf{x}_{i1} + \alpha_i)[1 - F(\boldsymbol{\beta}'\mathbf{x}_{i2} + \alpha_i)] + [1 - F(\boldsymbol{\beta}'\mathbf{x}_{i1} + \alpha_i)]F(\boldsymbol{\beta}'\mathbf{x}_{i2} + \alpha_i)$ . Because  $d_i$  is independently distributed with  $Ed_i = P_i$ , and both  $F$  and the variance of  $d_i$  are uniformly bounded, by a strong law of large numbers,

$$\frac{1}{N} J_{\tilde{B}_1} - \frac{1}{N} J \text{ almost surely } \rightarrow 0 \text{ as } N \rightarrow \infty \quad (6.3.20)$$

$$\text{if } \sum_{i=1}^N \frac{1}{i^2} \mathbf{m}_i \mathbf{m}_i' < \infty,$$

where  $\mathbf{m}_i$  replaces each element of  $(\mathbf{x}_{i2} - \mathbf{x}_{i1})$  by its square. The condition for convergence clearly holds if  $\mathbf{x}_{it}$  is uniformly bounded.

For the case of  $T > 2$ , there is no loss of generality in choosing the sequence  $D_{i1} = (d_{i11}, \dots, d_{i1T}, \sum_{t=1}^T d_{i1t} = \sum_{t=1}^T y_{it} = s, 1 \leq s \leq T-1)$  as the normalizing factor. Hence, we may rewrite the conditional probability (6.3.15) as

$$\text{Prob} \left( \mathbf{y}_i \mid \sum_{t=1}^T y_{it} \right) = \frac{\exp \{ \boldsymbol{\beta}' \sum_{t=1}^T \mathbf{x}_{it} (y_{it} - d_{i1t}) \}}{1 + \sum_{D_{ij} \in (\tilde{B}_i - D_{i1})} \exp \{ \boldsymbol{\beta}' \sum_{t=1}^T \mathbf{x}_{it} (d_{ijt} - d_{i1t}) \}} \quad (6.3.21)$$

Then the conditional log-likelihood function takes the form

$$\log L^* = \sum_{i \in C} \left\{ \boldsymbol{\beta}' \sum_{t=1}^T \mathbf{x}_{it} (y_{it} - d_{i1t}) - \log \left[ 1 + \sum_{D_{ij} \in (\tilde{B}_i - D_{i1})} \exp \left\{ \boldsymbol{\beta}' \sum_{t=1}^T \mathbf{x}_{it} (d_{ijt} - d_{i1t}) \right\} \right] \right\} \quad (6.3.22)$$

where  $C = \{i \mid \sum_{t=1}^T y_{it} \neq T, \sum_{t=1}^T y_{it} \neq 0\}$ .

Although we can find simple transformations of linear-probability and logit models that will satisfy the Neyman–Scott principle, we cannot find simple functions for the parameters



### 6.3.2 Random-Effects Models

When the individual specific effects  $\alpha_i$  are treated as random, we may still use the fixed-effects estimators to estimate the structural parameters  $\beta$ . The asymptotic properties of the fixed-effects estimators of  $\beta$  remain unchanged. However, if  $\alpha_i$  are random but are treated as fixed, the consequence at best is a loss of efficiency in estimating  $\beta$ . But it could be worse; namely, the resulting fixed-effects estimators may be inconsistent as discussed in Section 6.3.1.

When  $\alpha_i$  are independent of  $\mathbf{x}_i$  and are a random sampling from a univariate distribution  $G$ , indexed by a finite number of parameters  $\delta$ , the log-likelihood function becomes

$$\log L = \sum_{i=1}^N \log \int \prod_{t=1}^T F(\beta' \mathbf{x}_{it} + \alpha)^{y_{it}} [1 - F(\beta' \mathbf{x}_{it} + \alpha)]^{1-y_{it}} dG(\alpha | \delta). \quad (6.3.25)$$

where  $F(\cdot)$  is the distribution of the error term conditional on both  $\mathbf{x}_i$  and  $\alpha_i$ . Equation (6.3.25) replaces the probability function for  $y$  conditional on  $\alpha$  by a probability function that is marginal on  $\alpha$ . It is a function of a finite number of parameters  $(\beta', \delta')$ . Thus, maximizing (6.3.25), under weak regularity conditions, will give consistent estimators for  $\beta$  and  $\delta$  as  $N$  tends to infinity provided the distribution (or conditional distribution) of  $\alpha$  is correctly specified. If  $G(\alpha)$  is misspecified, maximizing (6.3.25) will yield inconsistent estimates when  $T$  is fixed. However, when  $T \rightarrow \infty$ , the random effects estimator becomes consistent, irrespective of the form of the postulated distribution of individual effects. The reason is that  $\alpha_i$  can practically be treated as a fixed constant for a given random draw,

$$\log f(y_i | \mathbf{x}_i, \beta, \alpha_i) = \sum_{t=1}^T \log f(y_{it} | \mathbf{x}_{it}, \beta, \alpha_i)$$

is a sum of  $T$  time series observation, so that the distribution of  $\alpha$  becomes negligible compared to that of the likelihood as the number of time periods increases (Arellano and Bonhomme 2009).

If  $\alpha_i$  is correlated with  $\mathbf{x}_{it}$ , maximizing (6.3.25) will not eliminate the omitted-variable bias. To allow for dependence between  $\alpha$  and  $\mathbf{x}$ , we must specify a distribution for  $\alpha$  conditional on  $\mathbf{x}$ ,  $G(\alpha | \mathbf{x})$  and consider the marginal log likelihood function

$$\log L = \sum_{i=1}^N \log \int \prod_{t=1}^T F(\beta' \mathbf{x}_{it} + \alpha)^{y_{it}} [1 - F(\beta' \mathbf{x}_{it} + \alpha)]^{1-y_{it}} dG(\alpha | \mathbf{x}) \quad (6.3.26)$$

A convenient specification suggested by Chamberlain (1980, 1984) is to assume that  $\alpha_i = \sum_{t=1}^T \mathbf{a}'_t \mathbf{x}_{it} + \eta_i = \mathbf{a}' \mathbf{x}_i + \eta_i$  where  $\mathbf{a}' = (\mathbf{a}'_1, \dots, \mathbf{a}'_T)$  and  $\mathbf{x}'_i = (\mathbf{x}'_{i1}, \dots, \mathbf{x}'_{iT})$ , and  $\eta_i$  is the residual. However, there is a very important difference in this step compared with the linear case. In the linear case it was not restrictive to decompose  $\alpha_i$  into its linear projection on  $\mathbf{x}_i$  and an orthogonal residual. Now we are assuming that the regression function  $E(\alpha_i | \mathbf{x}_i)$  is actually linear, that  $\eta_i$  is independent of  $\mathbf{x}_i$ , and that  $\eta_i$  has a specific probability distribution.

Given these assumptions, the log-likelihood function under our random-effects specification is

$$\log L = \sum_{i=1}^N \log \int \prod_{t=1}^T F(\beta' \mathbf{x}_{it} + \mathbf{a}' \mathbf{x}_i + \eta)^{y_{it}} \cdot [1 - F(\beta' \mathbf{x}_{it} + \mathbf{a}' \mathbf{x}_i + \eta)]^{1-y_{it}} dG^*(\eta), \quad (6.3.27)$$



$$\text{Prob}(y_{it} = 1) = \Phi \left[ (1 + \sigma_\eta^2)^{-1/2} (\boldsymbol{\beta}' \mathbf{x}_{it} + \mathbf{a}' \mathbf{x}_i) \right]. \quad (6.3.30)$$
$$\Pi = (1 + \sigma_\eta^2)^{-1/2} (I_T \otimes \boldsymbol{\beta}' + \mathbf{e}\mathbf{a}')$$

A more efficient estimator that avoids numerical integration is to impose the restriction (6.3.31) by  $\boldsymbol{\pi} = \text{vec}(\boldsymbol{\Pi}') = \mathbf{f}(\boldsymbol{\theta})$ , where  $\boldsymbol{\theta}' = (\boldsymbol{\beta}', \mathbf{a}', \sigma_\eta^2)$ , and use a GMM or minimum-distance estimator (see Chapter 3), just as in the linear case. Chamberlain (1984) suggests that we choose  $\hat{\boldsymbol{\theta}}$  to minimize<sup>13</sup>

$$(\hat{\boldsymbol{\pi}} - \mathbf{f}(\boldsymbol{\theta}))' \hat{\boldsymbol{\Omega}}^{-1} (\hat{\boldsymbol{\pi}} - \mathbf{f}(\boldsymbol{\theta})) \quad (6.3.32)$$

$$\Omega = J^{-1} \Delta J^{-1}, \quad (6.333)$$

$$J = \begin{bmatrix} J_1 & 0 & \dots & 0 \\ 0 & J_2 & & \\ \vdots & & \ddots & \\ 0 & & & J_T \end{bmatrix},$$

$$J_t = E \left\{ \frac{\phi_{it}^2}{\Phi_{it}(1 - \Phi_{it})} \mathbf{x}_i \mathbf{x}_i' \right\},$$

$$\Delta = E[\Phi_i \otimes \mathbf{x}_i \mathbf{x}_i'],$$

$$c_{it} = \frac{y_{it} - \Phi_{it}}{\Phi_{it}(1 - \Phi_{it})} \phi_{it}, \quad t = 1, \dots, T.$$

<sup>12</sup> In the case in which  $\alpha_i$  are uncorrelated with  $\mathbf{x}_i$ , we have  $\mathbf{a} = \mathbf{0}$  and  $\sigma_\eta^2 = \sigma_\alpha^2$ .

<sup>13</sup>  $\Omega$  is the asymptotic variance-covariance matrix of  $\hat{\pi}$  when no restrictions are imposed on the variance-covariance matrix of the  $T \times 1$  normal random variable  $\mathbf{u}_i + \epsilon \eta_i$ . We can relax the serial-independence assumption on  $u_{it}$  and allow  $E\mathbf{u}_i\mathbf{u}_i'$  to be an arbitrary positive definite matrix except for scale normalization. In this circumstance,  $\Pi = \text{diag}\{(\sigma_{\mu 1}^2 + \sigma_n^2)^{-1/2}, \dots, (\sigma_{\mu T}^2 + \sigma_n^2)^{-1/2}\}[I_T \otimes \beta' + e\alpha']$ .

## 6.4 SEMIPARAMETRIC APPROACH TO STATIC MODELS

The parametric approach of estimating discrete choice model suffers from two drawbacks: (1) conditional on  $\mathbf{x}$ , the probability law of generating  $(u_{it}, \alpha_i)$  is known a priori or conditional on  $\mathbf{x}$  and  $\alpha_i$ , the probability law of  $u_{it}$  is known a priori. (2) When  $\alpha_i$  are fixed, it appears that apart from logit and linear probability models, there is no simple transformation that can get rid of the incidental parameters. The semiparametric approach not only avoids making specific distribution of  $u_{it}$  but also allows a consistent estimator of  $\beta$  up to a scale whether  $\alpha_i$  is treated as fixed or random.

### 6.4.1 Maximum Score Estimator

Manski (1975, 1985, 1987) suggests a maximum score estimator that maximizes the sample average function

$$H_N(\mathbf{b}) = \frac{1}{N} \sum_{i=1}^N \sum_{t=2}^T \text{sgn}(\Delta \mathbf{x}'_{it} \mathbf{b}) \Delta y_{it} \quad (6.4.1)$$

subject to the normalization condition  $\mathbf{b}'\mathbf{b} = 1$ , where  $\Delta \mathbf{x}_{it} = \mathbf{x}_{it} - \mathbf{x}_{i,t-1}$ ,  $\Delta y_{it} = y_{it} - y_{i,t-1}$ ,  $\text{sgn}(w) = 1$  if  $w > 0$ , 0 if  $w = 0$ , and  $-1$  if  $w < 0$ . This is because under fairly general conditions (6.4.1) converges uniformly to

$$H(\mathbf{b}) = E[\text{sgn}(\Delta \mathbf{x}'_{it} \mathbf{b}) \Delta y_{it}], \quad (6.4.2)$$

where  $H(\mathbf{b})$  is maximized at  $\mathbf{b} = \beta^*$ , where  $\beta^* = \frac{\beta}{\|\beta\|}$  and  $\|\beta\|$  is the Euclidean norm  $\sum_{k=1}^K \beta_k^2$ .

To see this, we note that the binary choice model can be written in the form

$$y_{it} = \begin{cases} 1 & \text{if } y_{it}^* > 0, \\ 0 & \text{if } y_{it}^* \leq 0, \end{cases} \quad (6.4.3)$$

where  $y_{it}^*$  is given by (6.2.1) with  $v_{it} = \alpha_i + u_{it}$ . Under the assumption that  $u_{it}$  is independently identically distributed and is independent of  $\mathbf{x}_i$  and  $\alpha_i$  for given  $i$  (i.e.  $\mathbf{x}_{it}$  is strictly exogenous), we have

$$\begin{aligned} \mathbf{x}'_{it} \beta > \mathbf{x}'_{i,t-1} \beta &\iff E(y_{it} \mid \mathbf{x}_{it}) > E(y_{i,t-1} \mid \mathbf{x}_{i,t-1}), \\ \mathbf{x}'_{it} \beta = \mathbf{x}'_{i,t-1} \beta &\iff E(y_{it} \mid \mathbf{x}_{it}) = E(y_{i,t-1} \mid \mathbf{x}_{i,t-1}), \\ \mathbf{x}'_{it} \beta < \mathbf{x}'_{i,t-1} \beta &\iff E(y_{it} \mid \mathbf{x}_{it}) < E(y_{i,t-1} \mid \mathbf{x}_{i,t-1}). \end{aligned} \quad (6.4.4)$$

Rewrite (6.4.4) in terms of first differences, we have the equivalent representation

$$\begin{aligned} \Delta \mathbf{x}'_{it} \beta > 0 &\iff E[(y_{it} - y_{i,t-1}) > 0 \mid \Delta \mathbf{x}_{it}] \\ \Delta \mathbf{x}'_{it} \beta = 0 &\iff E[(y_{it} - y_{i,t-1}) = 0 \mid \Delta \mathbf{x}_{it}], \\ \Delta \mathbf{x}'_{it} \beta < 0 &\iff E[(y_{it} - y_{i,t-1}) < 0 \mid \Delta \mathbf{x}_{it}]. \end{aligned} \quad (6.4.5)$$

It is obvious that (6.4.5) continues to hold when  $\tilde{\beta} = \beta c$  where  $c > 0$ . Therefore, we shall only consider the normalized vector  $\beta^* = \frac{\beta}{\|\beta\|}$ .

Then, for any  $\mathbf{b}$  (satisfying  $\mathbf{b}'\mathbf{b} = 1$ ) such that  $\mathbf{b} \neq \beta^*$ ,

$$\begin{aligned} H(\beta^*) - H(\mathbf{b}) &= E[\{\text{sgn}(\Delta \mathbf{x}'_{it} \beta^*) - \text{sgn}(\Delta \mathbf{x}'_{it} \mathbf{b})\}(y_{it} - y_{i,t-1})] \\ &= 2 \int_{W_b} \text{sgn}(\Delta \mathbf{x}'_{it} \beta^*) E[y_{it} - y_{i,t-1} \mid \Delta \mathbf{x}] dF_{\Delta \mathbf{x}}, \end{aligned} \quad (6.4.6)$$



where  $W_b = [\Delta \mathbf{x} : \text{sgn}(\Delta \mathbf{x}' \boldsymbol{\beta}^*) \neq \text{sgn}(\Delta \mathbf{x}' \mathbf{b})]$ , and  $F_{\Delta \mathbf{x}}$  denotes the distribution of  $\Delta \mathbf{x}$ . Because of (6.4.5), the relation (6.4.6) implies that for all  $\Delta \mathbf{x}$ ,

$$\text{sgn}(\Delta \mathbf{x}' \boldsymbol{\beta}^*) E[y_t - y_{t-1} \mid \Delta \mathbf{x}] = |E[y_t - y_{t-1} \mid \Delta \mathbf{x}]|.$$

Therefore, under the assumption that  $\mathbf{x}$ 's are unbounded,<sup>14</sup>

$$H(\boldsymbol{\beta}^*) - H(\mathbf{b}) = 2 \int_{W_b} |E[y_t - y_{t-1} \mid \Delta \mathbf{x}]| dF_{\Delta \mathbf{x}} \geq 0. \quad (6.4.7)$$

Manski (1985, 1987) has shown that under fairly general conditions, the estimator maximizing the criterion function (6.4.1) is a strongly consistent estimator for  $\boldsymbol{\beta}^*$ .

As discussed in Chapter 3 and early sections of this chapter, when  $T$  is small, the MLE of the (structural) parameters  $\boldsymbol{\beta}$  is consistent as  $N \rightarrow \infty$  for the linear model and inconsistent for the nonlinear model in the presence of incidental parameters  $\alpha_i$ , because in the former case we can eliminate  $\alpha_i$  by differencing, while in the latter case we cannot. Thus, the error of estimating  $\alpha_i$  is transmitted into the estimator of  $\boldsymbol{\beta}$  in the nonlinear case. The Manski semiparametric approach makes use of the linear structure of the latent variable representation (6.2.1) or (6.4.4). The individual specific effects  $\alpha_i$  can again be eliminated by differencing and hence the lack of knowledge of  $\alpha_i$  no longer affects the estimation of  $\boldsymbol{\beta}$ .

The Manski maximum score estimator is consistent as  $N \rightarrow \infty$  for unknown conditional distribution of  $u_{it}$  given  $\alpha_i$  and  $\mathbf{x}_{it}, \mathbf{x}_{i,t-1}$ . However, it converges at the rate  $N^{1/3}$  which is much slower than the usual speed of  $N^{1/2}$  for the parametric approach. Moreover, Kim and Pollard (1990) have shown that  $N^{1/3}$  times the centered maximum score estimator converges in distribution to the random variable that maximizes a certain Gaussian process. This result shows that the maximum score estimator is probably not very useful in application since the properties of the limiting distribution are largely unknown.

The objective function (6.4.1) is equivalent to

$$\max_b H_N^*(\mathbf{b}) = N^{-1} \sum_{i=1}^N \sum_{t=2}^T [2 \cdot \mathbf{1}(\Delta y_{it} = 1) - 1] \mathbf{1}(\Delta \mathbf{x}_{it}' \mathbf{b} > 0) \quad (6.4.8)$$

subject to  $\mathbf{b}'\mathbf{b} = 1$ ,  $\mathbf{1}(A)$  is the indicator of the event  $A$  with  $\mathbf{1}(A) = 1$  if  $A$  occurs and 0 otherwise. The complexity of the maximum score estimator and its slow rate of convergence are due to the discontinuity of the function  $H_N(\mathbf{b})$  or  $H_N^*(\mathbf{b})$ . Horowitz (1992) suggests avoiding these difficulties by replacing  $H_N^*(\mathbf{b})$  with a sufficiently smooth function  $\tilde{H}_N(\mathbf{b})$  whose almost sure limit as  $N \rightarrow \infty$  is the same as that of  $H_N^*(\mathbf{b})$ . Let  $K(\cdot)$  be a continuous function of the real line  $(-\infty, \infty)$  into itself such that  $|K(v)| < M < \infty$  and  $\lim_{v \rightarrow -\infty} K(v) = 0$  and  $\lim_{v \rightarrow \infty} K(v) = 1$ . The  $K(\cdot)$  here is analogous to a cumulative distribution function. Let  $\{h_N : N = 1, 2, \dots\}$  be a sequence of strictly positive real numbers satisfying  $\lim_{N \rightarrow \infty} h_N = 0$ . Define

$$\tilde{H}_N(\mathbf{b}) = N^{-1} \sum_{i=1}^N \sum_{t=2}^T [2 \cdot \mathbf{1}(\Delta y_{it} = 1) - 1] K(\mathbf{b}' \Delta \mathbf{x}_{it} / h_N). \quad (6.4.9)$$

Horowitz (1992) defines a smoothed maximum score estimator as any solution that maximizes (6.4.9). Like Manski's estimator,  $\boldsymbol{\beta}$  can be identified only up to scale. Instead of using the normalization  $\|\boldsymbol{\beta}^*\| = 1$ , Horowitz (1992) finds it is more convenient to use the normalization that the coefficient of one component of  $\Delta \mathbf{x}$ , say  $\Delta x_1$ , to be equal to 1 in

<sup>14</sup> If  $\mathbf{x}$  is bounded, identification may fail if  $u_{it}$  is not logistic (Chamberlain 2010).

absolute value if its coefficient  $\beta_1 \neq 0$  and the probability distribution of  $\Delta \mathbf{x}_1$  conditional on the remaining components is absolutely continuous (with respect to Lebesgue measure).

The smoothed maximum score estimator is strongly consistent under the assumption that the distribution of  $\Delta u_{it} = u_{it} - u_{i,t-1}$  conditional on  $\Delta \mathbf{x}_{it}$  is symmetrically distributed with mean equal to zero. The asymptotic behavior of the estimator can be analyzed by taking a Taylor expansion of the first-order conditions and applying a version of the central limit theorem and the law of large numbers. The smoothed estimator of  $\boldsymbol{\beta}$  is consistent and, after centering and suitable normalization, is asymptotically normally distributed. Its rate of convergence is at least as fast as  $N^{-2/5}$  and, depending on how smooth the distribution of  $u$  and  $\boldsymbol{\beta}' \Delta \mathbf{x}$  are, can be arbitrarily close to  $N^{-1/2}$ .

#### 6.4.2 A Root-N Consistent Semiparametric Estimator

The speed of convergence of the smoothed maximum score estimator depends on the speed of convergence of  $h_N \rightarrow 0$ . Lee (1999) suggests a root- $N$  consistent semiparametric estimator that does not depend on a smoothing parameter by maximizing the double sums

$$\begin{aligned} & \{N(N-1)\}^{-1} \sum_{i \neq j} \sum_{t=2}^T \text{sgn}(\Delta \mathbf{x}'_{it} \mathbf{b} - \Delta \mathbf{x}'_{jt} \mathbf{b})(\Delta y_{it} - \Delta y_{jt}) \Delta y_{it}^2 \Delta y_{jt}^2 \\ &= \{N(N-1)\}^{-1} \sum_{\substack{i \\ i < j}} \sum_{\substack{j \\ \Delta y_{it} \neq \Delta y_{jt} \\ \Delta y_{it} \neq 0 \\ \Delta y_{jt} \neq 0}} \sum_{t=2}^T \text{sgn}(\Delta \mathbf{x}'_i \mathbf{b} - \Delta \mathbf{x}'_j \mathbf{b})(\Delta y_{it} - \Delta y_{jt}) \quad (6.4.10) \end{aligned}$$

with respect to  $\mathbf{b}$ . The consistency of the Lee estimator,  $\hat{\mathbf{b}}$ , follows from the fact that although  $\Delta y_i - \Delta y_j$  can take five values ( $0, \pm 1, \pm 2$ ), the event that  $(\Delta y_{it} - \Delta y_{jt}) \Delta y_{it}^2 \Delta y_{jt}^2 \neq 0$  excludes  $(0, \pm 1)$  to make  $\Delta y_{it} - \Delta y_{jt}$  binary (2 or  $-2$ ). Conditional on given  $j$ , the first average over  $i$  and  $t$  converges to

$$E\{\text{sgn}(\Delta \mathbf{x}' \mathbf{b} - \Delta \mathbf{x}'_j \mathbf{b})(\Delta y - \Delta y_j) \Delta y^2 \Delta y_j^2 \mid \Delta \mathbf{x}_j, \Delta y_j\} \quad (6.4.11)$$

The  $\sqrt{N}$  speed of convergence follows from the second average of the smooth function (6.4.10).

Normalizing  $\beta_1 = 1$ , the asymptotic covariance matrix of  $\sqrt{N}(\tilde{\mathbf{b}} - \tilde{\boldsymbol{\beta}})$  is equal to

$$4(E \nabla_2 \tau)^{-1} (E \nabla_1 \tau \nabla_1 \tau') (E \nabla_2 \tau)^{-1}, \quad (6.4.12)$$

where  $\tilde{\boldsymbol{\beta}} = (\beta_2, \dots, \beta_K)'$ , and  $\tilde{\mathbf{b}}$ , its estimator,

$$\tau(\Delta y_j, \Delta \mathbf{x}_j, \tilde{\mathbf{b}}) \equiv E_{i|j} \{\text{sgn}(\Delta \mathbf{x}'_i \mathbf{b} - \Delta \mathbf{x}'_j \mathbf{b})(\Delta y_i - \Delta y_j) \Delta y_i^2 \Delta y_j^2\}, \quad i \neq j,$$

with  $E_{i|j}$  denoting the conditional expectation of  $(\Delta y_i, \Delta \mathbf{x}'_i)$  conditional on  $(\Delta y_j, \Delta \mathbf{x}'_j)$ ,  $\nabla_1 \tau$  and  $\nabla_2 \tau$  denote the first and second derivative matrices of  $\tau(\Delta y_j, \Delta \mathbf{x}_j, \tilde{\mathbf{b}})$  with respect to  $\tilde{\mathbf{b}}$ .

The parametric approach requires the specification of the distribution of  $u$ . If the distribution of  $u$  is misspecified, the MLE of  $\boldsymbol{\beta}$  is inconsistent. The semiparametric approach does not require the specification of the distribution of  $u$  and permits its distribution to depend on  $\mathbf{x}$  in an unknown way (heteroskedasticity of unknown form). It is consistent up to a scale whether the unobserved individual effects are treated as fixed or correlated with  $\mathbf{x}$ . However, the step of differencing  $\mathbf{x}_{it}$  eliminates time-invariant variables from the estimation. Lee's (1999)  $\sqrt{N}$  consistent estimator takes the additional

differencing across individuals,  $\Delta \mathbf{x}_i - \Delta \mathbf{x}_j$ , further reduces the dimension of estimable parameters by eliminating “period individual-invariant” variables (e.g., time dummies and macroeconomic shocks common to all individuals) from the specification. Moreover, the requirement that  $u_{it}$  and  $u_{i,t-1}$  are identically distributed conditional on  $(\alpha_i, \mathbf{x}_{it}, \mathbf{x}_{i,t-1})$  does not allow the presence of the lagged dependent variables in  $\mathbf{x}_{it}$ . Neither can a semiparametric approach be used to generate the predicted probability conditional on  $\mathbf{x}$  as in the parametric approach. All it can estimate is the relative effects of the explanatory variables.

### 6.4.3 Estimation of the Coefficients of Time-Invariant Variables

The parametric or semiparametric estimates for binary choice (or single index) models with individual-specific effects are based on the difference of an individual's explanatory variables at two different periods. If the model contains time-invariant explanatory variables, the difference of an individual at explanatory variables at two time periods eliminates the time-invariant  $z_i$ . To obtain the coefficients of  $z_i, \gamma$ , requires a stronger assumption than those employed for linear regression models (see Section 2.6.1). Noting that, conditional on the relationship of the two different individual's time-varying variables being identical,  $\mathbf{x}'_{it}\boldsymbol{\beta} = \mathbf{x}'_{js}\boldsymbol{\beta}$ ,  $\alpha_i, \alpha_j, z_i, z_j$  are independently distributed, then similar to the argument (6.4.4),

$$P(y_{it} = 1 | \mathbf{x}'_{it} \boldsymbol{\beta} = a, \mathbf{z}_i) \geq P(y_{js} = 1 | \mathbf{x}'_{js} \boldsymbol{\beta} = a, \mathbf{z}_j) \quad (6.4.13)$$

if and only if

$$\mathbf{z}'_i \boldsymbol{\gamma} \geq \mathbf{z}'_j \boldsymbol{\gamma}. \quad (6.4.14)$$

Honoré and Kesina (2017) suggest to estimate  $\boldsymbol{\gamma}$  by

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \binom{N}{2}^{-1} \frac{1}{h_N} \sum_{i < j} \sum_{t=1}^{T_i} \sum_{s=1}^{T_j} K \left( \frac{(\mathbf{x}_{it} - \mathbf{x}_{js})' \hat{\boldsymbol{\beta}}}{h_N} \right) \cdot \text{sign}(\mathbf{y}_{it} - \mathbf{y}_{js}) \text{sign}[(\mathbf{z}_i - \mathbf{z}_j)' \mathbf{y}] \quad (6.4.15)$$

where  $K(\cdot)$  is a kernel function with  $|K(v)| < M$  for some constant  $M$  and  $K(v) \rightarrow 0$  as  $|v| \rightarrow \infty$ , and  $h_N$  is a bandwidth chosen so that in the limit, only pairs with  $\mathbf{x}'_{it}\boldsymbol{\beta} = \mathbf{x}'_{js}\boldsymbol{\beta}$  will contribute to the estimation of  $\hat{\boldsymbol{\gamma}}$ .

## 6.5 DYNAMIC MODELS

### 6.5.1 The General Model

The static models discussed in the previous sections assume that the probability of moving (or staying) in or out of a state is independent of the occurrence or nonoccurrence of the event in the past. However, in a variety of contexts, such as in the study of the incidence of accidents (Bates and Neyman 1951), brand loyalty (Chintagunta, Kyriazidou, and Perktold 2000), labor force participation (Heckman and Willis 1977; Hyslop 1999), and unemployment (Layton 1978), it is often noted that individuals who have experienced an event in the past are more likely to experience the event in the future than individuals who have not. In other words, the conditional probability that an individual will experience the event in the future is a function of past experience.

To analyze the intertemporal relationship among discrete variables, Heckman (1978, 1981b) proposed a general framework in terms of a latent-continuous-random-variable

crossing the threshold. He let the continuous random variable  $y_{it}^*$  be a function of  $\mathbf{x}_{it}$  and past occurrence of the event,

$$y_{it}^* = \beta' \mathbf{x}_{it} + \sum_{l=1}^{t-1} \gamma_l y_{i,t-l} + \phi \sum_{s=1}^{t-1} \prod_{l=1}^s y_{i,t-l} + v_{it}, \quad (6.5.1)$$

$$i = 1, \dots, N, t = 1, \dots, T$$

and

$$y_{it} = \begin{cases} 1 & \text{if } y_{it}^* > 0, \\ 0 & \text{if } y_{it}^* \leq 0. \end{cases} \quad (6.5.2)$$

The error term  $v_{it}$  is assumed to be independent of  $\mathbf{x}_{it}$  and is independently distributed over  $i$ , with a general intertemporal variance–covariance matrix  $E \mathbf{v}_i \mathbf{v}_i' = \Omega$ . The coefficient  $\gamma_l$  measures the effects of experience of the event  $l$  periods ago on current values of  $y_{it}^*$ . The coefficient  $\phi$  measures the effect of the cumulative recent spell of experience in the state for those still in the state on the current value of  $y_{it}^*$ .

Specifications (6.5.1) and (6.5.2) accommodate a wide variety of stochastic models that appear in the literature. For example, let  $\mathbf{x}_{it} = 1$ , and let  $v_{it}$  be independently identically distributed. If  $\gamma_l = 0, l = 2, \dots, T-1$ , and  $\phi = 0$ , Equations (6.5.1) and (6.5.2) generate a time-homogenous first-order Markov process. If  $\gamma_l = 0$  for all  $l$  and  $\phi \neq 0$ , a renewal process is generated. If  $\gamma_l = 0, l = 1, \dots, T-1$  and  $\phi = 0$ , a simple Bernoulli model results. If one allows  $v_{it}$  to follow an autoregressive moving-average scheme, but keeps the assumption that  $\gamma_l = 0, l = 1, \dots, T-1$ , and  $\phi = 0$ , the Coleman (1964) latent Markov model emerges.

As discussed before, repeated observations of a given group of individuals over time permit us to construct a model in which individuals may differ in their propensity to experience the event with the same source  $\mathbf{x}$ . Such heterogeneity is allowed by decomposing the error term  $v_{it}$  as

$$v_{it} = \alpha_i + u_{it}, \quad i = 1, \dots, N, t = 1, \dots, T, \quad (6.5.3)$$

where  $u_{it}$  is independently distributed over  $i$ , with arbitrary serial correlation, and  $\alpha_i$  is individual-specific and can be treated as a fixed constant or as random. Thus, for example, if the previous assumptions on the Markov process

$$\gamma_l = 0, l = 2, \dots, T-1, \text{ and } \phi = 0$$

hold, but  $v_{it}$  follows a “components-of-variance” scheme (6.5.3), a compound first-order Markov process, closely related to previous work on the mover-stayer model (Goodman 1961; Singer and Spilerman 1976), is generated.

Specifications (6.5.1)–(6.5.3) allow for three sources of persistence (after controlling for the observed explanatory variables,  $\mathbf{x}$ ). Persistence can be the result of serial correlation in the error term,  $u_{it}$ , or the result of “unobserved heterogeneity,”  $\alpha_i$ , or the result of true state dependence through the term  $\gamma_l y_{i,t-l}$  or  $\phi \prod_{l=1}^s y_{i,t-l}$ . Distinguishing the sources of persistence is important because a policy that temporarily increases the probability that  $y = 1$  will have different implications about future probabilities of experiencing an event.

When the conditional probability of an individual staying in a state is a function of past experience, two new issues arise. One is how to treat the initial observations. The second is how to distinguish true state dependence from spurious state dependence in which the past  $y_{it}$  appears in the specification merely as a proxy for the unobserved individual effects,  $\alpha_i$ . The first issue could play a role in deriving consistent estimators for a given model if

sample observations over time,  $T$ , are finite. The second issue is important because the time dependence among observed events could arise either from the fact that the actual experience of an event has modified individual behavior, or from unobserved components that are correlated over time, or from a combination of both.

### 6.5.2 Initial Conditions

When dependence among time-ordered outcomes is considered, just as in the dynamic linear regression model, the problem of initial conditions must be resolved for a likelihood approach before parameters generating the stochastic process can be estimated. In order to focus the discussion on the essential aspects of the problem, as in the previous example on the Markov process, where we assume  $\gamma_l = 0, l = 2, \dots, T - 1$ , and  $\phi = 0$  hold and no exogenous variables, so the observed data are generated by a first-order Markov process,

$$y_{it}^* = \beta_0 + \gamma y_{i,t-1} + v_{it}, \quad (6.5.4)$$

where

$$v_{it} = \alpha_i + u_{it},$$

$$y_{it} = \begin{cases} 1, & \text{if } y_{it}^* > 0, \\ 0, & \text{if } y_{it}^* \leq 0. \end{cases}$$

For ease of exposition, we shall also assume that  $u_{it}$  is independently normally distributed with mean zero and variance  $\sigma_u^2$  normalized to be equal to 1. It should be noted that the general conclusions of the following discussion also hold for other types of distributions.

In much applied work in the social sciences, two assumptions for initial conditions are typically invoked: (1) the initial conditions or relevant presample history of the process are assumed to be truly exogenous or (2) the process is assumed to be in equilibrium. Under the assumption that  $y_{i0}$  is a fixed nonstochastic constant for individual  $i$ , the joint probability of  $\mathbf{y}'_i = (y_{i1}, \dots, y_{iT})$ , given  $\alpha_i$ , is

$$\prod_{t=1}^T F(y_{it} \mid y_{i,t-1}, \alpha_i) = \prod_{t=1}^T \Phi\{(\beta_0 + \gamma y_{i,t-1} + \alpha_i)(2y_{it} - 1)\}, \quad (6.5.5)$$

where  $\Phi$  is the standard normal cumulative distribution function. Under the assumption that the process is in equilibrium, the limiting marginal probability for  $y_{it} = 1$  for all  $t$ , given  $\alpha_i$ , is (Karlin and Taylor 1975)<sup>15</sup>

<sup>15</sup> The transition-probability matrix of our homogeneous two-state Markov chain is

$$\mathcal{P} = \begin{bmatrix} 1 - \Phi(\beta_0 + \alpha_i) & \Phi(\beta_0 + \alpha_i) \\ 1 - \Phi(\beta_0 + \gamma + \alpha_i) & \Phi(\beta_0 + \gamma + \alpha_i) \end{bmatrix}.$$

By mathematical induction, the  $n$ -step transition matrix is

$$\mathcal{P}^n = \frac{1}{1 - \Phi(\beta_0 + \gamma + \alpha_i) + \Phi(\beta_0 + \alpha_i)} \times \left\{ \begin{bmatrix} 1 - \Phi(\beta_0 + \gamma + \alpha_i) & \Phi(\beta_0 + \alpha_i) \\ 1 - \Phi(\beta_0 + \gamma + \alpha_i) & \Phi(\beta_0 + \alpha_i) \end{bmatrix} + [\Phi(\beta_0 + \gamma + \alpha_i) - \Phi(\beta_0 + \alpha_i)]^n \times \begin{bmatrix} \Phi(\beta_0 + \alpha_i) & -\Phi(\beta_0 + \alpha_i) \\ -[1 - \Phi(\beta_0 + \gamma + \alpha_i)] & 1 - \Phi(\beta_0 + \gamma + \alpha_i) \end{bmatrix} \right\}.$$

$$P_i = \frac{\Phi(\beta_0 + \alpha_i)}{1 - \Phi(\beta_0 + \gamma + \alpha_i) + \Phi(\beta_0 + \alpha_i)}, \quad (6.5.6)$$

and the limiting probability for  $y_{it}=0$  is  $1 - P_i$ . Thus, the joint probability of  $(y_{i0}, \dots, y_{iT})$ , given  $\alpha_i$  is

$$\prod_{t=1}^T \Phi\{(\beta + \gamma y_{i,t-1} + \alpha_i)(2y_{it} - 1)\} P_i^{y_{i0}} (1 - P_i)^{1-y_{i0}}. \quad (6.5.7)$$

If  $\alpha_i$  is random, with distribution  $G(\alpha)$ , the likelihood function for the random-effects model under the first assumption is

$$L = \prod_{i=1}^N \int \prod_{t=1}^T \Phi\{(\beta_0 + \gamma y_{i,t-1} + \alpha)(2y_{it} - 1)\} dG(\alpha). \quad (6.5.8)$$

The likelihood function under the second assumption is

$$L = \prod_{i=1}^N \int \prod_{t=1}^T \Phi\{(\beta_0 + \gamma y_{i,t-1} + \alpha)(2y_{it} - 1)\} \cdot P_i^{y_{i0}} (1 - P_i)^{1-y_{i0}} dG(\alpha). \quad (6.5.9)$$

The likelihood functions (6.5.8) and (6.5.9) under both sets of assumptions about initial conditions are in closed form. When  $\alpha_i$  is treated as random, the MLEs for  $\beta_0, \gamma$ , and  $\sigma_\alpha^2$  are consistent if  $N$  tends to infinity or if both  $N$  and  $T$  tend to infinity. When  $\alpha_i$  is treated as a fixed constant (6.5.5), the MLEs for  $\beta_0, \gamma$ , and  $\alpha_i$  are consistent only when  $T$  tends to infinity. If  $T$  is finite, the MLE is biased. Moreover, the limited results from Monte Carlo experiments suggest that, contrary to the static case, the bias is significant (Heckman 1981b).

However, the assumption that initial conditions are fixed constants is justifiable only if the disturbances that generate the process are serially independent and if a genuinely new process is fortuitously observed at the beginning of the sample. If the process has been in operation prior to the time it is sampled, or if the disturbances of the model are serially dependent as in the presence of individual specific random effects, the initial conditions are not exogenous. The assumption that the process is in equilibrium also raises problems in many applications, especially when time-varying exogenous variables are driving the stochastic process.

Suppose that the analyst does not have access to the process from the beginning; then the initial state for individual  $i$ ,  $y_{i0}$ , can not be assumed fixed. The initial state is determined by the process generating the panel sample. The sample likelihood function for the fixed-effects model is

$$L = \prod_{i=1}^N \prod_{t=1}^T \Phi\{(\beta_0 + \gamma y_{i,t-1} + \alpha_i)(2y_{it} - 1)\} f(y_{i0} | \alpha_i), \quad (6.5.10)$$

and the sample likelihood function for the random-effects models is

$$L = \prod_{i=1}^N \int_{-\infty}^{\infty} \prod_{t=1}^T \Phi\{(\beta_0 + \gamma y_{i,t-1} + \alpha)(2y_{it} - 1)\} f(y_{i0} | \alpha) dG(\alpha), \quad (6.5.11)$$

where  $f(y_{i0} | \alpha)$  denotes the marginal probability of  $y_{i0}$  given  $\alpha_i$ . Thus, unless  $T$  is very large, maximizing (6.5.5) or (6.5.10) yields inconsistent estimates.<sup>16</sup>

Because  $y_{i0}$  is a function of unobserved past values, besides the fact that the marginal distribution of  $f(y_{i0} \mid \alpha)$  is not easy to derive, maximizing (6.5.10) or (6.5.11) is also considerably involved. Heckman (1981b) therefore suggested that we approximate the initial conditions for a dynamic discrete model by the following procedure:

*Step 1:* Approximate the probability of  $y_{i0}$ , the initial state in the sample, by a Probit model, with index function

$$y_{i0}^* = Q(\mathbf{x}_{it}) + \epsilon_{i0}, \quad (6.5.12)$$

and

$$y_{i0} = \begin{cases} 1 & \text{if } y_{i0}^* > 0, \\ 0 & \text{if } y_{i0}^* \leq 0, \end{cases} \quad (6.5.13)$$

where  $Q(\mathbf{x}_{it})$  is a general function of  $\mathbf{x}_{it}$ ,  $t = 0, \dots, T$ , usually specified as linear in  $\mathbf{x}_{it}$ , and  $\epsilon_{i0}$  is assumed to be normally distributed, with mean zero and variance 1.

*Step 2:* Permit  $\epsilon_{i0}$  to be freely correlated with  $v_{it}, t = 1, \dots, T$ .

*Step 3:* Estimate the model by maximum likelihood without imposing any restrictions between the parameters of the structural system and parameters of the approximate reduced-form probability for the initial state of the sample.

Heckman (1981b) conducted Monte Carlo studies comparing the performances of the MLEs when assumption on initial  $y_{i0}$  and  $\alpha_i$  conform with the true data generating process, an approximate reduced-form probability for  $y_{i0}$ , and false fixed  $y_{i0}$  and  $\alpha_i$  for a first-order Markov process. The data for his experiment were generated by the random-effects model

$$y_{it}^* = \beta x_{it} + \gamma y_{i,t-1} + \alpha_i + u_{it}, \quad (6.5.14)$$

$$y_{it} = \begin{cases} 1 & \text{if } y_{it}^* > 0, \\ 0 & \text{if } y_{it}^* \leq 0, \end{cases}$$

where the exogenous variable  $x_{it}$  was generated by (6.3.24). He let the process operate for 25 periods before selecting samples of 8 ( $= T$ ) periods for each of the 100 ( $= N$ ) individuals used in the 25 samples for each parameter set. Heckman's Monte Carlo results are produced in Table 6.2.

These results show that contrary to the static model, the fixed-effects Probit estimator performs poorly. The greater the variance of the individual effects ( $\sigma_\alpha^2$ ), the greater the bias. The  $t$  statistics based on the estimated information matrix also lead to a misleading inference by not rejecting the false null hypotheses of  $\gamma = \beta = 0$  in the vast majority of samples.

By comparison, Heckman's approximate solution performs better. Although the estimates are still biased from the true values, their biases are not significant, particularly when they are compared with the ideal estimates. The  $t$ -statistics based on the approximate solutions are also much more reliable than in the fixed-effects Probit model, because they lead to a correct inference in a greater proportion of the samples.

<sup>16</sup> This can be easily seen by noting that the expectation of the first-derivative vector of (6.5.5) or (6.5.8) with respect to the structural parameters does not vanish at the true parameter value when the expectations are evaluated under (6.5.10) or (6.5.11).

Table 6.2. Monte Carlo results for first-order Markov process

$\gamma$	$\sigma_{\alpha}^2 = 3$			$\sigma_{\alpha}^2 = 1$			
	$\beta = -0.1$	$\beta = 1$	$\beta = 0$	$\beta = -0.1$	$\beta = 1$	$\beta = 0$	
Values of $\hat{\gamma}$ and $\hat{\beta}$ for the random-effects estimator with known initial conditions <sup>a</sup>							
0.5	$\hat{\gamma}$	n.a. <sup>c</sup>	0.57	n.a. <sup>c</sup>			
	$\hat{\beta}$	n.a. <sup>c</sup>	0.94	— <sup>d</sup>			
0.1	$\hat{\gamma}$	0.13	0.12	0.14			
	$\hat{\beta}$	-0.11	1.10	—			
Values of $\hat{\gamma}$ and $\hat{\beta}$ for the approximate random-effects estimation <sup>a</sup>							
0.5	$\hat{\gamma}$	0.63	0.60	0.70	n.a. <sup>c</sup>	0.54	0.62
	$\hat{\beta}$	-0.131	0.91	—	n.a. <sup>c</sup>	0.93	—
0.1	$\hat{\gamma}$	0.14	0.13	0.17	0.11	0.11	0.13
	$\hat{\beta}$	-0.12	0.92	—	-0.12	0.95	—
Values of $\hat{\gamma}$ and $\hat{\beta}$ for the fixed-effects estimator <sup>b</sup>							
0.5	$\hat{\gamma}$	0.14	0.19	0.03	n.a. <sup>c</sup>	0.27	0.17
	$\hat{\beta}$	-0.07	1.21	—	n.a. <sup>c</sup>	1.17	—
0.1	$\hat{\gamma}$	-0.34	-0.21	-0.04	-0.28	-0.15	-0.01
	$\hat{\beta}$	-0.06	1.14	—	-0.08	1.12	—

<sup>a</sup>  $N = 100$ ;  $T = 3$ .

<sup>b</sup>  $N = 100$ ;  $T = 8$ .

<sup>c</sup> Data not available because the model was not estimated.

<sup>d</sup> Not estimated.

Source: Heckman (1981b, Table 4.2).

Heckman’s Monte Carlo results also point to a disquieting feature. The MLE produces a biased estimator even under the ideal conditions with a correctly specified likelihood function. Because a panel with 100 observations of three periods is not uncommon, this finding deserves further study.

### 6.5.3 A Conditional Approach

#### 6.5.3.1 Conditional Maximum Likelihood Estimator

The likelihood approach cannot yield a consistent estimator when  $T$  is fixed and  $N$  tends to infinity if the individual effects are fixed. If the individual effects are random and independent of  $\mathbf{x}$ , the consistency of the MLE depends on the correct formulation of the probability distributions of the effects and initial observations. A semiparametric approach cannot be implemented for a dynamic model because the strict exogeneity variables are violated with the presence of lagged dependent variables as explanatory variables. When the strict exogeneity condition of the explanatory variables is violated,  $E(\Delta u_{it} \mid \mathbf{x}_{it}, \mathbf{x}_{i,t-1}, y_{i,t-1}, y_{i,t-2}) \neq 0$ . In other words, the one-to-one correspondence relation of the form (6.4.4) is violated. Hence, the Manski (1985) type maximum score estimator cannot be implemented. Neither can the conditional approach be implemented. Consider the case of  $T = 2$ . The basic idea of the conditional approach is to consider the probability of  $y_{i2} = 1$  or 0 conditional on explanatory variables in both periods and conditional on  $y_{i1} \neq y_{i2}$ . If the explanatory variables of  $\text{Prob}(y_{i2} = 1)$  include  $y_{i1}$ , then the conditional



probability is either 1 or 0, accordingly as  $y_{i1} = 0$  or 1; hence it provides no information about  $\gamma$  and  $\beta$ .

However, in the case that  $T \geq 3$  and  $\mathbf{x}_{it}$  follows a certain special pattern, Honoré and Kyriazidou (2000) show that it is possible to generalize the conditional probability approach to consistently estimate the unknown parameters for the logit model or to generalize the maximum score approach without the need of formulating the distribution of  $\alpha_i$  or the probability distribution of the initial observations for certain types of discrete choice models. However, the estimators converge to the true values at the speed considerably slower than the usual square root  $N$  rate.

Consider the model (6.5.4) with the assumption that  $u_{it}$  is logistically distributed; then the model of  $(y_{i0}, \dots, y_{iT})$  is of the form

$$P(y_{i0} = 1 \mid \alpha_i) = P_0(\alpha_i) \quad (6.5.15)$$

$$P(y_{it} = 1 \mid \alpha_i, y_{i0}, \dots, y_{i,t-1}) = \frac{\exp(\gamma y_{i,t-1} + \alpha_i)}{1 + \exp(\gamma y_{i,t-1} + \alpha_i)}, \quad (6.5.16)$$

for  $t = 1, 2, \dots, T$ .

When  $T \geq 3$ , Chamberlain (1985) has shown that inference on  $\gamma$  can be made independent of  $\alpha_i$  by using a conditional approach.

For ease of exposition, we shall assume that  $T = 3$  (i.e., there are four time series observations for each  $i$ ). Consider the events

$$\begin{aligned} A &= \{y_{i0}, y_{i1} = 0, y_{i2} = 1, y_{i3}\}, \\ B &= \{y_{i0}, y_{i1} = 1, y_{i2} = 0, y_{i3}\}. \end{aligned}$$

where  $y_{i0}$  and  $y_{i3}$  can be either 1 or 0. Then

$$\begin{aligned} P(A) &= P_0(\alpha_i)^{y_{i0}} [1 - P_0(\alpha_i)]^{1-y_{i0}} \cdot \frac{1}{1 + \exp(\gamma y_{i0} + \alpha_i)} \\ &\quad \cdot \frac{\exp(\alpha_i)}{1 + \exp(\alpha_i)} \cdot \frac{\exp(y_{i3}(\gamma + \alpha_i))}{1 + \exp(\gamma + \alpha_i)} \end{aligned} \quad (6.5.17)$$

and

$$\begin{aligned} P(B) &= P_{i0}(\alpha_i)^{y_{i0}} [1 - P_0(\alpha_i)]^{1-y_{i0}} \\ &\quad \cdot \frac{\exp(\gamma y_{i0} + \alpha_i)}{1 + \exp(\gamma y_{i0} + \alpha_i)} \cdot \frac{1}{1 + \exp(\gamma + \alpha_i)} \cdot \frac{\exp(\alpha_i y_{i3})}{1 + \exp(\alpha_i)}. \end{aligned} \quad (6.5.18)$$

Hence

$$\begin{aligned} P(A \mid A \cup B) &= P(A \mid y_{i0}, y_{i1} + y_{i2} = 1, y_{i3}) \\ &= \frac{\exp(\gamma y_{i3})}{\exp(\gamma y_{i3}) + \exp(\gamma y_{i0})} \\ &= \frac{1}{1 + \exp[\gamma(y_{i0} - y_{i3})]}, \end{aligned} \quad (6.5.19)$$

and

$$\begin{aligned} P(B \mid A \cup B) &= P(B \mid y_{i0}, y_{i1} + y_{i2} = 1, y_{i3}) \\ &= 1 - P(A \mid A \cup B) \\ &= \frac{\exp[\gamma(y_{i0} - y_{i3})]}{1 + \exp[\gamma(y_{i0} - y_{i3})]}. \end{aligned} \quad (6.5.20)$$

Equations (6.5.19) and (6.5.20) are in the binary logit form and do not depend on  $\alpha_i$ . The conditional log-likelihood

$$\log \tilde{L} = \sum_{i=1}^N 1(y_{i1} + y_{i2} = 1) \{y_{i1} [\gamma(y_{i0} - y_{i3})] - \log [1 + \exp \gamma(y_{i0} - y_{i3})]\} \quad (6.5.21)$$

is in the conditional logit form. Maximizing (6.5.21) yields  $\sqrt{N}$  consistent estimator of  $\gamma$ , where  $1(A) = 1$  if  $A$  occurs and 0 otherwise.

When exogenous variables  $\mathbf{x}_{it}$  also appear as explanatory variables in the latent response function,

$$y_{it}^* = \boldsymbol{\beta}' \mathbf{x}_{it} + \gamma y_{i,t-1} + \alpha_i + u_{it}, \quad (6.5.22)$$

we may write

$$P(y_{i0} = 1 \mid \mathbf{x}_i, \alpha_i) = P_0(\mathbf{x}_i, \alpha_i), \quad (6.5.23)$$

$$\begin{aligned} P(y_{it} = 1 \mid \mathbf{x}_i, \alpha_i, y_{i0}, \dots, y_{i,t-1}) \\ = \frac{\exp(\mathbf{x}_{it}' \boldsymbol{\beta} + \gamma y_{i,t-1} + \alpha_i)}{1 + \exp(\mathbf{x}_{it}' \boldsymbol{\beta} + \gamma y_{i,t-1} + \alpha_i)}, \quad t = 1, \dots, T. \end{aligned} \quad (6.5.24)$$

Let  $P(y_{i0}) = P_0(\mathbf{x}_i, \alpha_i)^{y_{i0}} [1 - P_0(\mathbf{x}_i, \alpha_i)]^{1-y_{i0}}$ . Suppose  $T = 3$ . Under (6.5.24),

$$\begin{aligned} P(A) = P(y_{i0}) \cdot \frac{1}{1 + \exp(\mathbf{x}_{i1}' \boldsymbol{\beta} + \gamma y_{i0} + \alpha_i)} \cdot \frac{\exp(\mathbf{x}_{i2}' \boldsymbol{\beta} + \alpha_i)}{1 + \exp(\mathbf{x}_{i2}' \boldsymbol{\beta} + \alpha_i)} \\ \cdot \frac{\exp[(\mathbf{x}_{i3}' \boldsymbol{\beta} + \gamma + \alpha_i) y_{i3}]}{1 + \exp(\mathbf{x}_{i3}' \boldsymbol{\beta} + \gamma + \alpha_i)}. \end{aligned} \quad (6.5.25)$$

$$\begin{aligned} P(B) = P(y_{i0}) \cdot \frac{\exp(\mathbf{x}_{i1}' \boldsymbol{\beta} + \gamma y_{i0} + \alpha_i)}{1 + \exp(\mathbf{x}_{i1}' \boldsymbol{\beta} + \gamma y_{i0} + \alpha_i)} \\ \cdot \frac{1}{1 + \exp(\mathbf{x}_{i2}' \boldsymbol{\beta} + \gamma + \alpha_i)} \cdot \frac{[\exp(\mathbf{x}_{i3}' \boldsymbol{\beta} + \alpha_i) y_{i3}]}{1 + \exp(\mathbf{x}_{i3}' \boldsymbol{\beta} + \alpha_i)}. \end{aligned} \quad (6.5.26)$$

The denominators of  $P(A)$  and  $P(B)$  are different depending on whether the sequence is of  $(y_{i1} = 0, y_{i2} = 1)$  or  $(y_{i1} = 1, y_{i2} = 0)$ . Therefore, in general,  $P(A \mid \mathbf{x}_i, \alpha_i, A \cup B)$  will depend on  $\alpha_i$ . However, if  $\mathbf{x}_{i2} = \mathbf{x}_{i3}$ , then the denominators of  $P(A)$  and  $P(B)$  are identical. Using the same conditioning method, Honoré and Kyriazidou (2000) show that

$$\begin{aligned} P(A \mid \mathbf{x}_i, \alpha_i, A \cup B, \mathbf{x}_{i2} = \mathbf{x}_{i3}) \\ = \frac{1}{1 + \exp[(\mathbf{x}_{i1} - \mathbf{x}_{i2})' \boldsymbol{\beta} + \gamma(y_{i0} - y_{i3})]}, \end{aligned} \quad (6.5.27)$$

which does not depend on  $\alpha_i$ . If  $\mathbf{x}_{it}$  is continuous, it may be rare that  $\mathbf{x}_{i2} = \mathbf{x}_{i3}$ . Honoré and Kyriazidou (2000) propose estimating  $\boldsymbol{\beta}$  and  $\gamma$  by maximizing

$$\sum_{i=1}^N \mathbf{1}(y_{i1} + y_{i2} = 1) K \left( \frac{\mathbf{x}_{i2} - \mathbf{x}_{i3}}{h_N} \right) \ln \left\{ \frac{\exp[(\mathbf{x}_{i1} - \mathbf{x}_{i2})' \boldsymbol{\beta} + \gamma(y_{i0} - y_{i3})]^{y_{i1}}}{1 + \exp[(\mathbf{x}_{i1} - \mathbf{x}_{i2})' \boldsymbol{\beta} + \gamma(y_{i0} - y_{i3})]} \right\} \quad (6.5.28)$$

with respect to  $\beta$  and  $\gamma$  (over some compact set) if  $P(\mathbf{x}_{i2} = \mathbf{x}_{i3}) > 0$ . Here  $K(\cdot)$  is a kernel density function which gives appropriate weight to observation  $i$ , while  $h_N$  is a bandwidth which shrinks to zero as  $N$  tends to infinity. The asymptotic theory will require that  $K(\cdot)$  be chosen so that a number of regularity conditions are satisfied, such as  $|K(\cdot)| < M$  for some constant  $M$ , and  $K(v) \rightarrow 0$  as  $|v| \rightarrow \infty$  and  $\int K(v)dv = 1$ . For instance,  $K(v)$  is often taken to be the standard normal density function and  $h_N = cN^{-1/5}$  for some constant  $c$ . The effect of the term  $K\left(\frac{\mathbf{x}_{i2} - \mathbf{x}_{i3}}{h_N}\right)$  is to give more weight to observations for which  $\mathbf{x}_{i2}$  is close to  $\mathbf{x}_{i3}$ . Their estimator is consistent and asymptotically normal although their speed of convergence is only  $\sqrt{Nh_N^k}$  which is considerably slower than  $\sqrt{N}$  where  $k$  is the dimension of  $\mathbf{x}_{it}$ .

The conditional approach works for the logit model, but it does not seem applicable for general nonlinear models. However, if the nonlinearity can be put in the single index form  $F(a)$  with the transformation function  $F$  being a strictly increasing distribution function, then Manski's (1987) maximum score estimator for the static case can be generalized to the case where the lagged dependent variable is included in the explanatory variable set by considering

$$\begin{aligned} P(A \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3}) &= P_0(\mathbf{x}_i, \alpha_i)^{y_{i0}} [1 - P_0(\mathbf{x}_i, \alpha_i)]^{1-y_{i0}} \\ &\quad \times [1 - F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)] \times F(\mathbf{x}'_{i2}\beta + \alpha_i) \\ &\quad \times [1 - F(\mathbf{x}'_{i2}\beta + \gamma + \alpha_i)]^{1-y_{i3}} \times F(\mathbf{x}'_{i2}\beta + \gamma + \alpha_i)^{y_{i3}} \end{aligned} \quad (6.5.29)$$

and

$$\begin{aligned} P(B \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3}) &= P_0(\mathbf{x}_i, \alpha_i)^{y_{i0}} [1 - P_0(\mathbf{x}_i, \alpha_i)]^{1-y_{i0}} \\ &\quad \times F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i) \times [1 - F(\mathbf{x}'_{i2}\beta + \gamma + \alpha_i)] \\ &\quad \times [1 - F(\mathbf{x}'_{i2}\beta + \alpha_i)]^{1-y_{i3}} \times F(\mathbf{x}'_{i2}\beta + \alpha_i)^{y_{i3}}. \end{aligned} \quad (6.5.30)$$

If  $y_{i3} = 0$ , then

$$\begin{aligned} \frac{P(A \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3})}{P(B \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3})} &= \frac{[1 - F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)]}{[1 - F(\mathbf{x}'_{i2}\beta + \gamma + \alpha_i)]} \times \frac{F(\mathbf{x}'_{i2}\beta + \alpha_i)}{F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)} \\ &= \frac{[1 - F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)]}{[1 - F(\mathbf{x}'_{i2}\beta + \gamma y_{i3} + \alpha_i)]} \times \frac{F(\mathbf{x}'_{i2}\beta + \gamma y_{i3} + \alpha_i)}{F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)}, \end{aligned} \quad (6.5.31)$$

where the second equality follows from the fact that  $y_{i3} = 0$ . If  $y_{i3} = 1$ , then

$$\begin{aligned} \frac{P(A \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3})}{P(B \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3})} &= \frac{[1 - F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)]}{[1 - F(\mathbf{x}'_{i2}\beta + \gamma + \alpha_i)]} \times \frac{F(\mathbf{x}'_{i2}\beta + \gamma + \alpha_i)}{F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)} \\ &= \frac{[1 - F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)]}{[1 - F(\mathbf{x}'_{i2}\beta + \gamma y_{i3} + \alpha_i)]} \times \frac{F(\mathbf{x}'_{i2}\beta + \gamma y_{i3} + \alpha_i)}{F(\mathbf{x}'_{i1}\beta + \gamma y_{i0} + \alpha_i)}, \end{aligned} \quad (6.5.32)$$

where the second equality follows from the fact that  $y_{i3} = 1$ , so that  $\gamma y_{i3} = \gamma$ . In either case, the monotonicity of  $F$  implies that

$$\frac{P(A)}{P(B)} \begin{cases} > 1 \text{ if } \mathbf{x}'_{i2}\boldsymbol{\beta} + \gamma y_{i3} > \mathbf{x}'_{i1}\boldsymbol{\beta} + \gamma y_{i0}, \\ < 1 \text{ if } \mathbf{x}'_{i2}\boldsymbol{\beta} + \gamma y_{i3} < \mathbf{x}'_{i1}\boldsymbol{\beta} + \gamma y_{i0}. \end{cases}$$

Therefore,

$$\begin{aligned} & \text{sgn}[P(A \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3}) - P(B \mid \mathbf{x}_i, \alpha_i, \mathbf{x}_{i2} = \mathbf{x}_{i3})] \\ &= \text{sgn}[(\mathbf{x}_{i2} - \mathbf{x}_{i1})'\boldsymbol{\beta} + \gamma(y_{i3} - y_{i0})]. \end{aligned} \quad (6.5.33)$$

Hence, Honoré and Kyriazidou (2000) propose a maximum score estimator that maximizes the score function

$$\sum_{i=1}^N K\left(\frac{\mathbf{x}_{i2} - \mathbf{x}_{i3}}{h_N}\right) (y_{i2} - y_{i1}) \text{sgn}[(\mathbf{x}_{i2} - \mathbf{x}_{i1})'\boldsymbol{\beta} + \gamma(y_{i3} - y_{i0})] \quad (6.5.34)$$

with respect to  $\boldsymbol{\beta}$  and  $\gamma$ . Honoré and Kyriazidou's estimator is consistent (up to a scale) if the density of  $\mathbf{x}_{i2} - \mathbf{x}_{i3}$ ,  $f(\mathbf{x}_{i2} - \mathbf{x}_{i3})$ , is strictly positive at zero,  $f(0) > 0$ . (This assumption is required for consistency.)

We have discussed the estimation of panel data dynamic discrete choice model assuming that  $T = 3$ . It can be easily generalized to the case of  $T > 3$  by maximizing the objective function that is based on sequences where an individual switches between alternatives in any two of the middle  $T - 1$  periods:

$$\begin{aligned} & \sum_{i=1}^N \sum_{1 \leq s < t \leq T-1} \mathbf{1}\{y_{is} + y_{it} = 1\} K\left(\frac{\mathbf{x}_{i,t+1} - \mathbf{x}_{i,s+1}}{h_N}\right) \\ & \times \ln \left( \frac{\exp[(\mathbf{x}_{is} - \mathbf{x}_{it})'\boldsymbol{\beta} + \gamma(y_{i,s-1} - y_{i,t+1}) + \gamma(y_{i,s+1} - y_{i,t-1})\mathbf{1}(t-s > 1)]^{y_{is}}}{1 + \exp[(\mathbf{x}_{is} - \mathbf{x}_{it})'\boldsymbol{\beta} + \gamma(y_{i,s-1} - y_{i,t+1}) + \gamma(y_{i,s+1} - y_{i,t-1})\mathbf{1}(t-s > 1)]} \right) \end{aligned} \quad (6.5.35)$$

The conditional approach does not require modeling of the initial observations of the sample. Neither does it make any assumptions about the statistical relationship of the individual effects with the observed explanatory variables or with the initial conditions. However, it also suffers from the limitation that  $\mathbf{x}_{is} - \mathbf{x}_{it}$  has support in a neighborhood of 0 for any  $t \neq s$  which rules out time-dummies as explanatory variables.<sup>17</sup> The fact that individual effects cannot be estimated also means that it is not possible to carry out predictions or to compute elasticities for individual agents at specified values of the explanatory variables or to evaluate the impact of a social policy.

### 6.5.3.2 A Cohort Approach to Estimate the Individual-Specific Effects

The Honoré and Kyriazidou (2000) (HK) conditional maximum likelihood estimator allows an investigator to obtain consistent estimators of the coefficients of the dynamic quality choice models where the time-varying conditional variables include a lag dependent variable when  $N$  is large and  $T$  is finite. The HK approach eliminated  $\alpha_i$  in the conditional likelihood function, hence it cannot estimate the unobserved individual-specific effects,  $\alpha_i$  or coefficients of time-invariant variables. Damrongplasit and Hsiao (2021) use a cohort approach proposed by Deaton (1985) in which all individuals belonging to the  $g - th$  cohort have identical  $\alpha_i = \alpha_g$  to estimate  $\alpha_g$  in the logit model.

<sup>17</sup> See Arellano and Carrasco (2003) for a GMM approach to estimate the dynamic random-effects Probit model.

The logit model implies that

$$\log \frac{p_{it}}{1 - p_{it}} = \gamma y_{i,t-1} + \beta' \mathbf{x}_{it} + \alpha_i, \quad (6.5.36)$$

However, to implement the logistic regression of (6.5.36), the probability can neither be 0 nor 1. Damrongplasit and Hsiao (2021) divide observations into  $G$  subgroups with the assumption that for those individuals belonging to the  $g$ th subgroup have identical  $\alpha_i = \alpha_g$ . They approximate the logit model by

$$\log \frac{p_{gt}}{1 - p_{gt}} = \lambda \bar{y}_{g,t-1} + \beta' \bar{\mathbf{x}}_{gt} + \alpha_g + \varepsilon_{gt}, \quad g = 1, \dots, G, \quad (6.5.37)$$

where

$$\bar{y}_{g,t-1} = \frac{1}{N_{gt}} \sum_{i=1}^N 1(y_{it} \in g) y_{i,t-1}, \quad \bar{\mathbf{x}}_{gt} = \frac{1}{N_{gt}} \sum_{i=1}^N 1(\mathbf{x}_{it} \in g) \mathbf{x}_{it},$$

and  $N_{gt} = \sum_{i=1}^N 1(y_{it} \in g)$ , where  $1(\cdot)$  denotes the indicator function. Estimate  $p_{gt}$  by

$$\hat{p}_{gt} = \frac{1}{N_{gt}} \sum_{i=1}^N 1(y_{it} \in g) y_{it}, \quad g = 1, \dots, G. \quad (6.5.38)$$

Substituting  $\hat{p}_{gt}$  in lieu of  $p_{gt}$  in (6.5.37), the cohort-specific effects,  $\alpha_g$ , are estimated by

$$\min \sum_{t=1}^T \sum_{g=1}^G \left[ \log \frac{\hat{p}_{gt}}{1 - \hat{p}_{gt}} - \hat{\gamma} \bar{y}_{g,t-1} - \hat{\beta}' \bar{\mathbf{x}}_{gt} - \alpha_g \right]^2. \quad (6.5.39)$$

Minimizing (6.5.39), conditional on  $\gamma$  and  $\beta$ , yields the least squares estimate of  $\alpha_g$

$$\hat{\alpha}_g = \frac{1}{GT} \sum_{t=1}^T \sum_{g=1}^G \left( \log \frac{\hat{p}_{gt}}{1 - \hat{p}_{gt}} - \hat{\gamma} \bar{y}_{g,t-1} - \hat{\beta}' \bar{\mathbf{x}}_{gt} \right). \quad (6.5.40)$$

#### 6.5.4 State dependence versus heterogeneity

There could be three possible scenarios for observing that  $\text{Prob}(y_{it} = 1)$  given  $y_{i,t-1} = 1$  is different from  $\text{Prob}(y_{it} = 1)$  given  $y_{i,t-1} = 0$  and different from the marginal  $\text{Prob}(y_{it} = 1)$ . The first is that the probability of observing  $y_{it}$  is altered due to experiencing an event leading to changes the preferences, pricing, constraints or parameters relevant for future choice (Lucas 1976).<sup>18</sup> The second is that the probability of observing  $y_{it}$  depends on past state  $y_{i,t-1}$  or more periods before. The third is that individuals may differ in certain unmeasured variables that influence their probability of experiencing the event but are not influenced by the experience of the event. If these variables are correlated over time and are not properly controlled, previous experience may appear to be a determinant of future experience solely because it is a proxy for such temporally persistent unobservables. We shall refer to the first scenario as a *structural break* or a *true state dependence*, i.e. whether the parameters characterizing  $\text{Prob}(y_{it})$  stay constant or have changed due to either  $y_{i,t-1} = 0$  or 1. The second scenario is termed as the *true dynamic dependence* and the third scenario *spurious dynamic dependence*, because in the former case, past experience

<sup>18</sup> In this sense, either a static model or a dynamic model could be subject to a structural break.

has a genuine behavioral effect in the sense that an otherwise identical individual who has experienced the event will behave differently in the next period or following periods than an individual who has not experienced the event. In the latter case, previous experience appears to be a determinant of future experience solely because it is a proxy for temporally persistent unobservables that determine choices (e.g., Heckman 1978, 1981a, 1981c).

The problem of distinguishing between a *structural break* and the *true or spurious state dependencies* is of considerable substantive interest. To demonstrate this, let us consider some work in the theory of unemployment. Phelps (1972) argued that current unemployment has a real and lasting effect on the probability of future unemployment. Hence, short-term economic policies that alleviate unemployment tend to lower aggregate unemployment rates in the long run by preventing the loss of work-enhancing market experience. On the other hand, Cripps and Tarling (1974) maintained the opposite view in their analysis of the incidence and duration of unemployment. They assumed that individuals differ in their propensity to experience unemployment and in their unemployment duration times, and those differences can not be fully accounted for by measured variables. They further assumed that the actual experience of having been unemployed or the duration of past unemployment does not effect future incident or duration (in the long run). Hence, in their model, short-term economic policies have no effect on long-term unemployment. Underlying the controversy between Phelps (1972) and Cripps and Tarling (1974) are two issues: (i) whether experiencing an event leads to changes in the behavioral parameters; and (ii) if there is a dynamic dependence in the sense that the probability of occurrence depends on whether an individual experienced an event in the last period or the periods before.

We shall call whether experiencing an event leads to changes in the behavioral parameters a *structural break*. If whether experiencing an event changes the probability of experiencing the event in the next period or beyond independent of whether the behavioral parameters stay constant or not, the *true dynamic dependence*, i.e., the coefficients of lag dependent variables, differs from zero. If experiencing an event leads to changes in the behavioral parameters, the long run or equilibrium conditions are changed independent of whether a model is static or dynamic. If experiencing an event does not lead to changes in behavioral parameters and if the dynamic dependence is stationary, it leads only to changes in the short-run probability. It does not alter the equilibrium condition.

The empirical observation that  $\text{Prob}(y_{it}) \neq \text{Prob}(y_{it}|y_{i,t-1})$  could be due either to a *structural break* (when experiencing an event leads to changes in behavioral parameters),<sup>19</sup> or to *true dynamic dependence* (when the parameters of past state occurrence are different from zero) or to *spurious dynamic dependence* (where coefficients of lag dependent variables are different from zero because of the presence of unobserved individual effects,  $\alpha_i$  that persist over time). Ignoring these effects of unmeasured variables (heterogeneity) creates serially correlated residuals. This suggests that we cannot use the conditional probability, given past occurrence not equal to the marginal probability alone,  $\text{Prob}(y_{it} | y_{i,t-s}, \mathbf{x}_{it}) \neq \text{Prob}(y_{it} | \mathbf{x}_{it})$ , to test for true state dependence against spurious state dependence. A proper test for true versus spurious state dependence should control for the time-persistent unobserved individual-specific effects. So is a test for structural break of parameter homogeneity before and after experiencing an event.

When conditional on the individual effects,  $\alpha_i$ , the error term  $u_{it}$  is serially uncorrelated, a test for dynamic dependence can be implemented by controlling the individual effects and testing for the conditional probability equal to the marginal probability,

<sup>19</sup> Note that a structural break could occur in a static model too.

$$\text{Prob}(y_{it} \mid y_{i,t-s}, \mathbf{x}_{it}, \boldsymbol{\theta}, \alpha_i) = \text{Prob}(y_{it} \mid \mathbf{x}_{it}, \boldsymbol{\theta}, \alpha_i). \quad (6.5.41)$$

where  $\boldsymbol{\theta}$  denotes the parameter vector of a model. When  $N$  is fixed and  $T \rightarrow \infty$ , likelihood ratio tests can be implemented to test (6.5.41).<sup>20</sup> However, if  $T$  is finite, controlling  $\alpha_i$  to obtain a consistent estimator for the coefficient of a lagged dependent variable imposes very restrictive conditions on the data which could severely limit the power of the test, as shown in Section 6.4.

If  $\alpha_i$  are treated as random and the conditional distribution of  $\alpha_i$  given  $\mathbf{x}_i$  is known, a more powerful test is to use an unconditional approach. Thus, one may test true dynamic dependence versus spurious dynamic dependence by testing the significance of the MLE of  $\gamma$  of the log likelihood

$$\sum_{i=1}^N \log \int \prod_{t=1}^T \left\{ F(\mathbf{x}'_{it}\boldsymbol{\beta} + \gamma y_{i,t-1} + \alpha_i)^{y_{it}} \left[ 1 - F(\mathbf{x}'_{it}\boldsymbol{\beta} + \gamma y_{i,t-1} + \alpha_i) \right]^{1-y_{it}} \cdot P(\mathbf{x}_i, \alpha)^{y_{i0}} \left[ 1 - P(\mathbf{x}_i, \alpha) \right]^{1-y_{i0}} \right\} G(\alpha_i \mid \mathbf{x}_i) d\alpha_i. \quad (6.5.42)$$

When conditional on the individual effects,  $\alpha_i$ , the error term  $u_{it}$  remains serially correlated, the problem becomes more complicated. The conditional probability,  $\text{Prob}(y_{it} \mid y_{i,t-1}, \alpha_i)$ , not being equal to the marginal probability,  $\text{Prob}(y_{it} \mid \alpha_i)$ , could be because of past  $y_{it}$  containing information on  $u_{it}$ . A test for dynamic dependence cannot simply rely on the multinomial distribution of the  $(y_{i1}, \dots, y_{iT})$  sequence. The general framework (6.5.1) and (6.5.2) proposed by Heckman (1978, 1981a, 1981b) accommodates very general sorts of heterogeneity and structural dependence. It permits an analyst to combine models and test among competing specifications within a unified framework. However, the computations of maximum likelihood methods for the general models could be quite involved. It would be useful to rely on simple methods to explore data before implementing the computationally cumbersome maximum-likelihood method for a specific model.

Chamberlain (1978) suggested a simple method to distinguish true dynamic dependence from spurious dynamic dependence. He noted that just as in the continuous models, a key distinction between dynamic dependence and serial correlation is whether or not there is a dynamic response to an intervention. This distinction can be made clear by examining (6.5.1). If  $\gamma = 0$ , a change in  $\mathbf{x}$  has its full effect immediately, where as if  $\gamma \neq 0$ , this implies a distributed-lag response to a change in  $\mathbf{x}$ . The lag structure relating  $y$  to  $\mathbf{x}$  is not related to the serial correlation in  $u$ . If  $\mathbf{x}$  is increased in period  $t$  and then returned to its former level, the probability of  $y_{i,t+1}$  is not affected if  $\gamma = 0$ , because by assumption the distribution of  $u_{it}$  was not affected. If  $\gamma \neq 0$ , then the one-period shift in  $\mathbf{x}$  will have lasting effects. An intervention that affects the probability of  $y$  in period  $t$  will continue to affect the probability of  $y$  in period  $t+1$ , even though the intervention was presented only in period  $t$ . In contrast, an interpretation of serial correlation is that the shocks ( $u$ ) tend to persist for more than one period and that  $y_{i,t-s}$  is informative only in helping to infer  $u_{it}$

<sup>20</sup> Let  $P_{it} = \text{Prob}(y_{it} \mid \mathbf{x}_{it}, \alpha_i)$  and  $P_{it}^* = \text{Prob}(y_{it} \mid y_{i,t-\ell}, \mathbf{x}_{it}, \alpha_i)$ . Let  $\hat{P}_{it}$  and  $\hat{P}_{it}^*$  be the MLEs obtained by maximizing  $\mathcal{L} = \prod_i \prod_t P_{it}^{y_{it}} (1 - P_{it})^{1-y_{it}}$  and  $\mathcal{L}^* = \prod_i \prod_t P_{it}^{*y_{it}} (1 - P_{it}^*)^{1-y_{it}}$  with respect to unknown parameters, respectively. A likelihood-ratio test statistic for the null hypothesis (6.5.41) is  $-2 \log [\mathcal{L}(\hat{P}_{it}) / \mathcal{L}(\hat{P}_{it}^*)]$ . When conditional on  $\mathbf{x}_{it}$  and  $\alpha_i$ , there are repeated observations; we can also use the Pearson chi-square goodness-of-fit statistic to test (6.5.41). For details, see Bishop, Fienberg, and Holland (1975, chapter 7). However, in the finite- $T$  case, the testing procedure cannot be implemented, as the  $\alpha_i$ 's are unknown and cannot be consistently estimated.

Table 6.3. *Estimates of employment models for women aged 45–59 in 1968<sup>a</sup>*

Variable	(1)	(2)	(3)
Intercept	–2.576 (4.6)	1.653 (2.5)	0.227 (0.4)
No. of children aged <6	–0.816 (2.7)	–0.840 (2.3)	–0.814 (2.1)
County unemployment rate (%)	–0.035 (1.5)	–0.027 (1.0)	–0.018 (0.57)
County wage rate (\$/h)	0.104 (0.91)	0.104 (0.91)	0.004 (0.02)
Total no. of children	–0.146 (4.3)	–0.117 (2.2)	–0.090 (2.4)
Wife’s education (years)	0.162 (6.5)	0.105 (2.8)	0.104 (3.7)
Family income, excluding wife’s earnings	$-0.363 \times 10^{-4}$ (4.8)	$-0.267 \times 10^{-4}$ (2.7)	$-0.32 \times 10^{-4}$ (3.6)
National unemployment rate	–0.106 (0.51)	–0.254 (1.4)	–1.30 (6)
Recent experience	0.143 (0.95)	0.273 (1.5)	1.46 (12.2)
Predicted presample experience	0.072 (5.8)	0.059 (3.4)	0.045 (3.4)
Serial-correlation coefficient:			
$\rho_{12}$	0.913	—	—
$\rho_{13}$	0.845		
$\rho_{23}$	0.910		
$\rho$	—	0.873 (14.0)	—
$\sigma_\alpha^2/(\sigma_u^2 + \sigma_\alpha^2)$	—	—	—
Log likelihood	–237.74	–240.32	–263.65

<sup>a</sup> Asymptotic normal test statistics in parentheses; these statistics were obtained from the estimating information matrix.

and hence to predict  $u_{it}$ . Therefore, a test that is not very sensitive to functional form is to simply include lagged  $\mathbf{x}'$ s without lagged  $y$ . After conditioning on the individual-specific effect  $\alpha$ , there may be two possible outcomes. If there is no dynamic dependence, then

$$\text{Prob}(y_{it} = 1 \mid \mathbf{x}_{it}, \mathbf{x}_{i,t-1}, \dots, \alpha_i) = \text{Prob}(y_{it} = 1 \mid \mathbf{x}_{it}, \alpha_i). \tag{6.5.43}$$

If there is dynamic dependence, then

$$\text{Prob}(y_{it} = 1 \mid \mathbf{x}_{it}, \mathbf{x}_{i,t-1}, \dots, \alpha_i) \neq \text{Prob}(y_{it} = 1 \mid \mathbf{x}_{it}, \alpha_i). \tag{6.5.44}$$

While the combination of (6.5.43) and (6.5.44) provides a simple form to distinguish pure heterogeneity, dynamic dependence, and serial correlation, we can not make further distinctions with regard to different forms of dynamic dependence, heterogeneity, and serial correlation should (6.5.43) be rejected. A structural break test has to be conducted to test Phelps’s (1972) contention that experiencing an event leads to change in the behavioral parameters.

6.5.5 Three Examples

The control of heterogeneity plays a crucial role in distinguishing true dynamic dependence from spurious dynamic dependence. Neglecting heterogeneity and the issue of initial observations can also seriously bias the coefficient estimates. It is important in estimating dynamic models that the heterogeneity in the sample be treated correctly. To demonstrate this, we use the female-employment models estimated by Heckman (1981c), household brand choices estimated by Chintagunta, Kyriazidou, and Perktold (2001), and labor participation decision and health shock analysis by Damrongplasit and Hsiao (2021).



Table 6.3. (cont.)

(4)	(5)	(6)	(7)	(8)
-2.367 (6.4)	-2.011 (3.4)	-2.37 (5.5)	-3.53 (4.6)	-1.5 (0)
-0.742 (2.6)	-0.793 (2.1)	-0.70 (2.0)	-1.42 (2.3)	-0.69 (1.2)
-0.030 (1.5)	-0.027 (1.2)	-0.03 (1.6)	-0.059 (1.3)	0.046 (11)
0.090 (0.93)	0.139 (1.5)	0.13 (1.4)	0.27 (1.1)	0.105 (0.68)
-0.124 (4.9)	-0.116 (2.2)	-0.161 (4.9)	-0.203 (3.9)	-0.160 (6.1)
0.152 (7.3)	0.095 (2.5)	0.077 (3)	0.196 (4.8)	0.105 (3.3)
$-0.312 \times 10^{-4}$ (5.2)	$-0.207 \times 10^{-4}$ (2.3)	$-0.2 \times 10^{-4}$ (2.6)	$-0.65 \times 10^{-4}$ (5.1)	$-0.385 \times 10^{-4}$ (20)
-0.003 (0.38)	-0.021 (0.26)	0.02 (3)	1.03 (0.14)	-0.71 (0)
— <sup>b</sup>	—	—	—	—
0.062 (0.38)	0.062 (3.5)	0.091 (7.0)	0.101 (5.4)	0.095 (11.0)
0.917	—	—	—	—
0.873	—	—	—	—
0.946	—	—	—	—
—	-0.942 (50)	—	—	—
—	—	0.92 (4.5)	—	0.941 (4.1)
-239.81	-243.11	-244.7	-367.3	-242.37

<sup>b</sup> Not estimated.

Source: Heckman (1981c, Table 3.2).

### 6.5.5.1 Female Employment

Heckman (1981c) used the first three-year sample of women aged 45–59 in 1968 from the Michigan Panel Survey of Income dynamics to study married women's employment decisions. A woman is defined to be a market participant if she worked for money any time in the sample year. The set of explanatory variables is as follows: the woman's education; family income, excluding the wife's earnings; number of children younger than six; number of children at home; unemployment rate in the county in which the woman resided; the wage of unskilled labor in the county (a measure of the availability of substitutes for a woman's time in the home); the national unemployment rate for prime-age males (a measure of aggregate labor-market tightness); two types of prior work experience, within-sample work experience and presample work experience. The effect of previous work experience is broken into two components, because it is likely that presample experience exerts a weaker measured effect on current participation decisions than more recent experience. Furthermore, because the data on presample work experience are based on a retrospective question and therefore are likely to be measured with error, Heckman replaces them with predicted values based on a set of regressors.

Under the assumption that behavioral parameters stay constant over time, Heckman fit the data to various multivariate Probit models of the form (6.5.1) and (6.5.2) to investigate whether or not work experience raises the probability that a woman will work in the (short-run) future (by raising her wage rates) and to investigate the importance of controlling for heterogeneity in utilizing panel data. Maximum-likelihood-coefficient estimates for the state-dependent models under the assumptions of stationary intertemporal covariance matrix

$$\Omega = \begin{bmatrix} 1 & \rho_{12} & \rho_{13} \\ & 1 & \rho_{23} \\ & & 1 \end{bmatrix},$$

first-order Markov process ( $v_{it} = \rho v_{i,t-1} + u_{it}$ ), and no heterogeneity ( $v_{it} = u_{it}$ ) are presented in columns 1, 2, and 3, respectively, of Table 6.3.<sup>21</sup> Coefficient estimates for no state dependence with general stationary intertemporal correlation, first-order Markov process, conventional random effects error-component formulation  $v_{it} = \alpha_i + u_{it}$ , equivalent to imposing the restriction that  $\rho_{12} = \rho_{13} = \rho_{23} = \sigma_\alpha^2 / (\sigma_u^2 + \sigma_\alpha^2)$ , and no heterogeneity are presented in columns 4, 5, 6, and 7, respectively. A Heckman–Willis (1977) model with time-invariant exogenous variables and conventional error-component formulation was also estimated and is presented in column 8.

Likelihood ratio test statistics (twice the difference of the log-likelihood value) against the most general model (column 1, Table 6.3) indicate the acceptance of recent labor-market experience as an important determinant of the current employment decision, with unobservables determining employment choices following a first-order Markov process (column 2, Table 6.3) as a favored hypothesis, and the statistics clearly reject all other formulations. In other words, Heckman's study found that work experience, as a form of general and specific human capital investment, raises the probability that a woman will more likely work in the near future, even after accounting for serial correlation of a very general type. It also maintained that unobserved variables exist that affect labor participation.

Comparison of the estimates of the favored hypothesis with estimates of other models indicates that the effect of recent market experience on employment is dramatically overstated in a model that neglects heterogeneity. The estimated effect of recent market experience on current employment status recorded in column 3, Table 6.3, overstates the impact by a factor of 10 (1.46 vs. 0.143)! Too much credit will be attributed to past experience as a determinant of employment if intertemporal correlation in the unobservables is ignored. Likewise for the estimated impact of national unemployment on employment. On the other hand, the effect of children on employment is understated in models that ignore heterogeneity.

Comparisons of various models' predictive performance on sample-run patterns (temporal employment status) are presented in Table 6.4. It shows that dynamic models ignoring heterogeneity underpredict the number of individuals who work all of the time and overpredict the number who do not work at all. It also overstates the estimated frequency of turnover in the labor force. In fact, comparing the performances of the predicted run patterns for the dynamic and static models without heterogeneity (column 3 and 7 or Table 6.3 and columns 3 and 4 of Table 6.4) suggests that introducing "lagged employment status" into a model as a substitute for a more careful treatment of heterogeneity is an imperfect procedure. In this case, it is worse than using no proxy at all. Nor does a simple static model with a "components-of-variance" scheme (column 8 of Table 6.3, column 5 of Table 6.4) perform any better. Dynamic models that neglect heterogeneity (column 3 of Table 6.4) overestimate labor-market turnover, whereas the static model with a conventional variance components formulation (column 5 of Table 6.4) overstates the extent of heterogeneity and the degree of intertemporal correlation. It overpredicts the number who never work during these three years and underpredicts the number who always work.

21 A nonstationary model was also estimated by Heckman (1981c), but because the data did not reject stationarity, we shall treat the model as having stationary covariance.







Table 6.6. Maximum-likelihood estimators of dynamic model with random effects

	Pooled sample, ALC = 1	Pooled sample, ALC = 0	Male, ALC = 1	Male, ALC = 0	Female, ALC = 1	Female, ALC = 0
Lagged labor force participation	3.354 [0.056]***	2.835 [0.052]***	3.653 [0.082]***	2.839 [0.092]***	3.054 [0.074]***	2.611 [0.064]***
Male	0.643 [0.053]***	0.957 [0.053]***				
Age 55–59	–0.872 [0.073]***	–0.851 [0.092]***	–0.999 [0.110]***	–1.013 [0.162]***	–0.852 [0.098]***	–0.907 [0.114]***
Age 60–64	–1.908 [0.086]***	–1.766 [0.104]***	–2.003 [0.123]***	–2.087 [0.168]***	–1.913 [0.121]***	–1.753 [0.139]***
Age 65 above	–3.295 [0.108]***	–3.462 [0.127]***	–3.205 [0.153]***	–3.908 [0.201]***	–3.472 [0.152]***	–3.394 [0.176]***
Married	–0.078 [0.054]	0.321 [0.058]***	0.263 [0.090]***	1.107 [0.120]***	–0.333 [0.071]***	–0.127 [0.068]*
Year12	0.448 [0.078]***	0.768 [0.067]***	0.342 [0.125]***	1.017 [0.116]***	0.591 [0.103]***	0.709 [0.085]***
Postschool	0.606 [0.061]***	0.937 [0.068]***	0.346 [0.085]***	1.044 [0.119]***	0.828 [0.088]***	0.866 [0.085]***
Degree	0.974 [0.079]***	1.282 [0.077]***	0.702 [0.118]***	1.288 [0.145]***	1.195 [0.107]***	1.321 [0.094]***
Younger child	–0.783 [0.075]***	–1.343 [0.064]***	–0.167 [0.135]	–0.27 [0.165]	–1.23 [0.099]***	–1.702 [0.079]***
Older child	–0.012 [0.059]	0.212 [0.057]***	0.132 [0.099]	0.445 [0.147]***	–0.201 [0.077]***	0.086 [0.065]
ln(income)	0.57 [0.032]***	0.315 [0.031]***	0.576 [0.047]***	0.305 [0.051]***	0.589 [0.045]***	0.359 [0.040]***
Unemployment rate	–0.016 [0.019]	–0.013 [0.019]	–0.012 [0.029]	0.004 [0.035]	–0.018 [0.026]	–0.014 [0.024]
Health shock	–0.792 [0.058]***	–0.173 [0.104]*	–0.944 [0.084]***	–0.351 [0.166]**	–0.689 [0.081]***	–0.072 [0.136]
Constant	–7.054 [0.361]***	–4.046 [0.360]***	–6.809 [0.522]***	–3.543 [0.596]***	–6.923 [0.500]***	–3.934 [0.459]***
Number of observations	42,285	42,309	19,782	19,637	22,503	22,672

(1) Standard errors are in parentheses.

(2) \*\*\* Significant at 1%. \*\* Significant at 5%. \* Significant at 10%.

(3) ALC = 1 implies that there is a long-term health condition, while ALC = 0 implies that there is no long-term health condition.

Source: Damrongplasit and Hsiao (2021, Table 2).

rejected the homogeneity between two groups. The Hausman (1978) specification test between the random-effects and fixed-effects specification also rejected the random-effects specification.

Figures 6.1 and 6.2 reproduce the plot by Damrongplasit and Hsiao (2021) for the dynamic response path to a health shock for the male group, aged 55–59, married, with post-school education, with younger children having the mean of the  $\ln(\text{income})$  and the mean of unemployment rate at period  $t$  based on random effects and fixed effects estimated coefficients where the cohort specific effects are estimated using (6.5.40) and the recursive formula is suggested by Damrongplasit et al. (2019),

$$P(y_{it} = 1) = P(y_{i,t-1} = 1)P(y_{it} = 1 \mid y_{i,t-1} = 1, \mathbf{x}'_{it}) \\ + [1 - P(y_{i,t-1} = 1)]P(y_{it} = 1 \mid y_{i,t-1} = 0, \mathbf{x}'_{it}). \quad (6.5.50)$$

These figures show the dynamic response path, although it depends on the initial condition (whether labor force participation decision in the initial period is 0 or 1,  $LFP(0) = 0$

Table 6.7. *Maximum-likelihood estimators of dynamic model with fixed effects (bandwidth parameter = 8)*

	Pooled sample, ALC = 1	Pooled sample, ALC = 0	Male, ALC = 1	Male, ALC = 0	Female, ALC = 1	Female, ALC = 0
Lagged labor force participation	1.705 [0.104]***	1.65 [0.124]***	1.707 [0.175]***	1.518 [0.248]***	1.707 [0.135]***	1.716 [0.150]***
Age 55–59	–0.621 [0.347]*	–1.193 [0.567]**	–0.456 [0.549]	–2.427 [1.479]	–0.738 [0.459]	–0.73 [0.649]
Age 60–64	–1.549 [0.517]***	–2.104 [0.814]***	–1.68 [0.854]**	–3.827 [1.739]**	–1.449 [0.662]**	–1.298 [1.019]
Age 65 above	–2.38 [0.664]***	–3.288 [1.060]***	–2.528 [1.011]**	–5.29 [2.031]***	–2.214 [0.957]**	–2.154 [1.368]
Married	0.01 [0.261]	–0.332 [0.334]	0.041 [0.475]	0.353 [0.670]	–0.059 [0.318]	–0.618 [0.405]
Year12	0.762 [0.402]*	1.361 [0.326]***	1.902 [0.742]**	1.711 [0.565]***	0.018 [0.519]	1.048 [0.418]**
Postschool	1.126 [0.420]***	1.154 [0.473]**	0.845 [0.767]	0.868 [0.917]	1.161 [0.522]**	1.302 [0.584]**
Degree	1.277 [0.685]*	2.3 [0.654]***	2.794 [1.273]**	2.045 [0.943]**	0.132 [0.881]	3.027 [1.147]***
Younger child	–1.054 [0.260]***	–1.399 [0.241]***	–0.163 [0.511]	–0.303 [0.745]	–1.395 [0.319]***	–1.513 [0.266]***
Older child	–0.119 [0.255]	–0.149 [0.290]	0.37 [0.500]	–1.148 [0.979]	–0.342 [0.305]	–0.058 [0.310]
ln(income)	0.176 [0.112]	0.314 [0.144]**	0.175 [0.180]	0.235 [0.250]	0.178 [0.147]	0.355 [0.182]*
Unemployment rate	0.106 [0.061]*	–0.085 [0.072]	0.192 [0.103]*	–0.032 [0.148]	0.053 [0.077]	–0.083 [0.085]
Health shock	–0.621 [0.154]***	–0.204 [0.315]	–0.679 [0.236]***	–0.727 [0.602]	–0.608 [0.212]***	0.034 [0.399]
Number of observations	5391	4270	2176	1237	3215	3033

(1) Standard errors are in parentheses.

(2) \*\*\* Significant at 1%. \*\* Significant at 5%. \* Significant at 10%.

(3) ALC = 1 implies that there is a long-term health condition, while ALC = 0 implies that there is no long-term health condition.

Source: Damrongplasit and Hsiao (2021, Table 3a).

or 1) and whether an individual suffers a health shock or not in the initial period,  $(HS(0) = 1 \text{ or } 0)$ , the long-run equilibrium conditions remain unchanged. However, there is a substantial difference between the group not suffering activity limiting condition and the group that do. There is also a substantial difference between the random-effects and fixed-effects estimates. In other words, their exercise shows that it is important to take account of both the observed and unobserved sample heterogeneity to construct an econometric model that fits the data. It also shows that whether a policy change (or external change) can alter the long-run equilibrium is an issue of parameter stability in light of a policy change.

## 6.6 ALTERNATIVE APPROACHES FOR IDENTIFYING DYNAMIC DEPENDENCE

Section 6.5 focuses on getting consistent estimators for dynamic panel discrete choice models with individual specific effects. If individual specific effects are treated as random, the consistency of dynamic models requires the knowledge of the conditional distribution of individual specific effects  $\alpha_i$  given the  $T$  time series observations of the  $K \times 1$  exogenous variables,  $\mathbf{x}_{it}, \mathbf{x}'_i = (\mathbf{x}'_{i1}, \dots, \mathbf{x}'_{iT})$ ,  $G(\alpha_i | \mathbf{x}_i)$ , and the initial value distribution

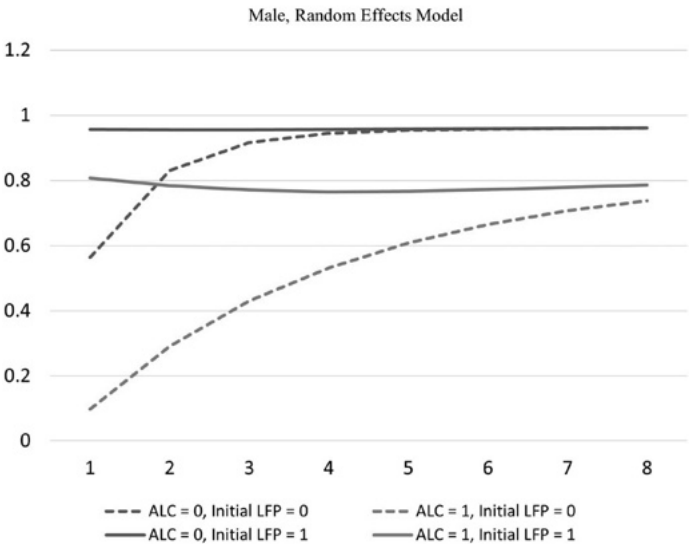


Figure 6.1. Male dynamic response path for the health shock that leads to activity limiting (ALC = 1) vs. not (ALC = 0).

- (1) Initial LFP = 0 implies that an individual was out of the labor force in the initial period, while Initial LFP = 1 implies that an individual was in the labor force in the initial period.
- (2) ALC = 0 implies that there is no long-term health condition (i.e., activity limiting condition), while ALC = 1 implies that there is a long-term health condition.

Source: Damrongplasit and Hsiao (2021, Figure 1).

given  $\mathbf{x}_i$ ,  $P(y_{i0} | \mathbf{x}_i)$ . If  $\alpha_i$  is treated as a fixed constant, the consistency of the MLE requires  $T \rightarrow \infty$ . If  $T$  is finite, the conditions for obtaining a consistent estimator of the coefficients of exogenous variables and lagged dependent variables impose severe restrictions on the observed data that only a very small proportion of the sample may be utilized, if they satisfy the conditions at all. In this section, we consider alternative approaches to identify the dynamic dependence–bias reduced estimator for fixed-effects models; bounding parameters without the knowledge of  $G(\alpha | \mathbf{x})$  and  $P(y_0 | \mathbf{x})$  for random-effects models; and approximate model.

6.6.1 Bias Adjusted Estimator

Controlling the impact of unobserved heterogeneity in linear models is relatively straightforward (e.g., see Chapters 2 and 3). Controlling the impact of unobserved heterogeneity that is correlated with explanatory variables in nonlinear models is much more difficult. When  $T$  is finite, the estimators of the parameters of interest (structural parameters) are inconsistent no matter how large  $N$  is. This inconsistency occurs because only a finite number of observations are available to estimate each individual effect  $\alpha_i$ , while the estimation of structural parameters depends on  $\alpha_i$ . Increasing  $T$  does not necessarily fully solve this problem if  $N$  also grows with  $T$  (e.g., see Alvarez and Arellano 2003; Hahn and Kuersteiner 2011; Hahn and Newey 2004). In this section we consider methods that reduce the bias of the estimator to the order of  $\frac{1}{T^2}$ .



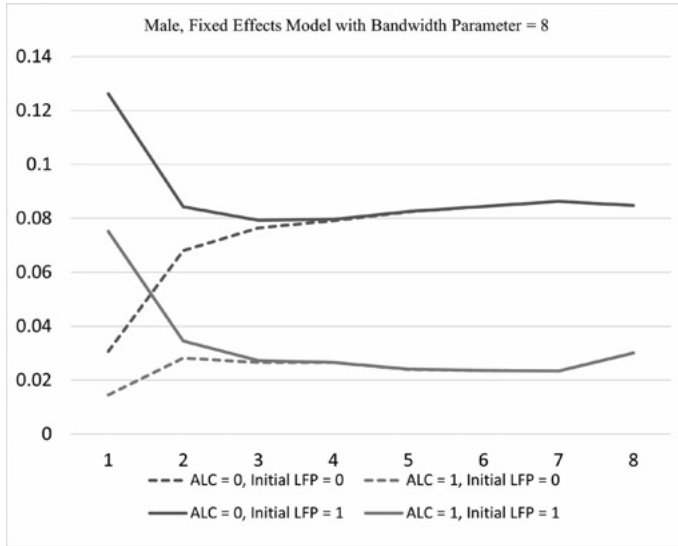


Figure 6.2. Male dynamic response path for the health shock that leads to activity limiting (ALC = 1) vs. not (ALC = 0).

- (1) Initial LFP = 0 implies that an individual was out of the labor force in the initial period, while Initial LFP = 1 implies that an individual was in the labor force in the initial period.
- (2) ALC = 0 implies that there is no long-term health condition (i.e., activity limiting condition), while ALC = 1 implies that there is a long-term health condition.

Source: Damrongplasit and Hsiao (2021, Figure 2).

Let  $\theta$  denote the parameters of interest (structural parameters) and  $\alpha_i$  denote the unobserved individual-specific effects (incidental parameters). Let  $\hat{\theta}_T$  denote the estimator of  $\theta$  based on  $NT$  panel data  $(y_{it}, x_{it})$  and  $\hat{\alpha}_i$ , the estimated  $\alpha_i$ , say the fixed effects MLE for the static logit model (6.3.3 and 6.3.4) or the dynamic logit model (6.5.16). In general, because of the error in the estimation of  $\alpha_i$  when  $T$  is fixed, under the assumption that cross-sectional units are independent over  $i$  and  $\hat{\theta}_T \rightarrow \theta_T$  when  $N \rightarrow \infty$ , where

$$\theta_T = \theta + \frac{B}{T} + \frac{D}{T^2} + O\left(\frac{1}{T^3}\right) \quad (6.6.1)$$

for some  $B$  and  $D$ . This bias could be small for large  $T$ . However, if  $N$  grows at the same rate as  $T$  when  $T \rightarrow \infty$ , the fixed-effects estimator  $\hat{\theta}$  is asymptotically biased. For  $\frac{N}{T} \rightarrow c \neq 0$ ,

$$\sqrt{NT}(\hat{\theta} - \theta) = \sqrt{NT}(\hat{\theta} - \theta_T) + \sqrt{NT} \cdot \frac{B}{T} + O\left(\sqrt{\frac{N}{T^3}}\right) \quad (6.6.2)$$

will have asymptotic normal distribution centered at  $\sqrt{c}B$ . (e.g., the fixed-effects estimator for the dynamic panel data model (3.2.11)).

Hahn and Newey (2004) suggest a delete one-panel Jackknife estimator to reduce the bias,

$$\tilde{\theta} \equiv T\hat{\theta}_T - \frac{T-1}{T} \sum_{t=1}^T \hat{\theta}(t), \quad (6.6.3)$$

where  $\hat{\theta}(t)$  is the fixed-effects estimator based on the subsample excluding the observations of the  $t$ th period. If  $\theta_T$  has the form (6.6.1), then the estimator  $\hat{\theta}$  will converge in probability to

$$\begin{aligned} & (T\theta_T - (T-1)\theta_{T-1}) \\ &= \theta + \left(\frac{1}{T} - \frac{1}{T-1}\right)D + O\left(\frac{1}{T^2}\right) \\ &= \theta + O\left(\frac{1}{T^2}\right). \end{aligned} \quad (6.6.4)$$

Thus, the Jackknife estimator reduces the bias to the order of  $\frac{1}{T^2}$ . However, in addition to the fact that the Jackknife estimator (6.6.3) requires the estimation of  $(T+1)$  fixed-effects estimators, the asymptotic covariance matrix of  $\hat{\theta}$  is complicated to derive unless  $(y_{it}, \mathbf{x}_{it})$  are contemporaneously and intertemporally independently distributed (over  $i$  and  $t$ ). To allow the dependence over  $t$  without complicating the derivation of a symptotic variance of Jackknife bias corrected estimator, Dhaene and Jochmans (2015) suggest using subpanels formed by *consecutive* observations for each cross-sectional unit, the *split-panel Jackknife estimation*. Their idea is to split the complete panel with  $T$  time series observations into  $S$  nonoverlapping subpanels,  $S \geq 2$ , where the element of each subpanel is *consecutive* to preserve the time dependence structure of the full panel. Let  $\hat{\theta}_s$  denote the subpanel estimate of  $s$ th block; then

$$E\left(\frac{T_s}{T}\hat{\theta}_s - \hat{\theta}_T\right) = \frac{B}{T} + o\left(\frac{1}{T}\right), \quad s = 1, \dots, S, \quad (6.6.5)$$

where  $T_s$  denotes the number of time series dimensions in block  $s$ . Then

$$\tilde{\theta} = S\hat{\theta}_T - (S-1)\left(\frac{1}{S}\sum_{s=1}^S\left(\frac{T_s}{T}\right)\hat{\theta}_s\right) \quad (6.6.6)$$

reduces the bias to the smaller order of  $\left(\frac{1}{T}\right), o\left(\frac{1}{T}\right)$ .

An alternative approach is to obtain an estimated  $B, \hat{B}$ , then forming a bias-corrected estimator

$$\hat{\theta}^* = \hat{\theta} - \frac{\hat{B}}{T}, \quad (6.6.7)$$

(e.g., (3.8.27)). The advantage of (6.6.7) is that it reduces the bias, but the formula for the asymptotic covariance matrix of  $\hat{\theta}^*$  remains the same as that of  $\hat{\theta}$ . However, the derivation of  $\hat{B}$  can be complicated.

For instance, consider the panel dynamic binary choice model of the form,

$$\begin{aligned} y_{it} &= 1(\mathbf{x}_{it}'\beta + y_{i,t-1}\gamma + \alpha_i + u_{it} > 0), \\ i &= 1, \dots, N, \quad t = 1, \dots, T, \quad y_{i0} \text{ observable}, \end{aligned} \quad (6.6.8)$$

where  $1(A) = 1$  if event  $A$  occurs and 0 otherwise. We suppose that  $u_{it}$  is independently, identically distributed with mean 0. Then

$$\begin{aligned} E(y_{it} \mid y_{i,t-1}, \mathbf{x}_{it}, \alpha_i) &= \text{Prob}(y_{it} = 1 \mid y_{i,t-1}, \mathbf{x}_{it}, \alpha_i) \\ &= F(\mathbf{x}_{it}'\beta + y_{i,t-1}\gamma + \alpha_i) \\ &= F_{it}, \end{aligned} \quad (6.6.9)$$

where  $F$  is the integral of the probability distribution function of  $u_{it}$  from  $-(\mathbf{x}_{it}\boldsymbol{\beta} + y_{i,t-1}\gamma + \alpha_i)$  to  $\infty$ . When  $\alpha_i$  is considered fixed, the log-likelihood function conditional on  $y_{i0}$  takes the form

$$\log L = \sum_{i=1}^N \sum_{t=1}^T [y_{it} \log F_{it} + (1 - y_{it}) \log (1 - F_{it})] \quad (6.6.10)$$

The MLE of  $\boldsymbol{\theta}' = (\boldsymbol{\beta}', \gamma)$  and  $\alpha_i$  are obtained by solving the following first-order conditions simultaneously:

$$\frac{\partial \log L}{\partial \alpha_i} \big|_{\hat{\alpha}_i} = 0, \quad i = 1, \dots, N, \quad (6.6.11)$$

$$\frac{\partial \log L}{\partial \boldsymbol{\theta}} \big|_{\hat{\boldsymbol{\theta}}} = \mathbf{0}. \quad (6.6.12)$$

Substituting the solutions of (6.6.11) as a function of  $\boldsymbol{\theta}$  to (6.6.10) yields the concentrated log-likelihood function

$$\log L^* = \sum_{i=1}^N \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})), \quad (6.6.13)$$

where

$$\begin{aligned} \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})) = \sum_{t=1}^T \bigg\{ & y_{it} \log F(\mathbf{x}'_{it}\boldsymbol{\beta} + y_{i,t-1}\gamma + \hat{\alpha}_i(\boldsymbol{\theta})) \\ & + (1 - y_{it}) \log [1 - F(\mathbf{x}'_{it}\boldsymbol{\beta} + y_{i,t-1}\gamma + \hat{\alpha}_i(\boldsymbol{\theta}))] \bigg\}. \end{aligned}$$

Then the MLE of  $\boldsymbol{\theta}$  is the solution of the following first-order conditions:

$$\frac{1}{NT} \sum_{i=1}^N \left[ \frac{\partial \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} + \frac{\partial \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \hat{\alpha}_i(\boldsymbol{\theta})} \times \frac{\partial \hat{\alpha}_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right]_{\hat{\boldsymbol{\theta}}} = \mathbf{0}. \quad (6.6.14)$$

The estimating equation (6.6.14) depends on  $\hat{\alpha}_i$ . When  $T \rightarrow \infty$ ,  $\hat{\alpha}_i \rightarrow \alpha_i$ , the MLE of  $\boldsymbol{\theta}$  is consistent. When  $T$  is finite,  $\hat{\alpha}_i \neq \alpha_i$ , then (6.6.14) evaluated at  $\boldsymbol{\theta}$  does not converge to zero. Hence, the MLE of  $\boldsymbol{\theta}$ ,  $\hat{\boldsymbol{\theta}}$ , is not consistent. The bias of the MLE is of order  $\frac{1}{T}$ . The analytical solution for  $\hat{\boldsymbol{\theta}}_T$  can be derived by taking a Taylor series expansion of (6.6.14) (e.g. see Hahn and Kuersteiner 2011).

Instead of obtaining  $\hat{\boldsymbol{\theta}}_T$  directly, Carro (2007) proposes to derive the bias-corrected MLE directly by taking the Taylor series expansion of the score function (6.6.14) around  $\alpha_i$  and evaluating it at the true value  $\boldsymbol{\theta}$ , which yields

$$\begin{aligned} d_{\theta_i}(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})) &= d_{\theta_i}(\boldsymbol{\theta}, \alpha_i) + d_{\theta\alpha_i}(\hat{\alpha}_i(\boldsymbol{\theta}) - \alpha_i) \\ &+ \frac{1}{2} d_{\theta\alpha_i\alpha_i}(\hat{\alpha}_i(\boldsymbol{\theta}) - \alpha_i)^2 + O_p(T^{-\frac{1}{2}}), i = 1, \dots, N. \end{aligned} \quad (6.6.15)$$

where  $d_{\theta_i}(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})) = \frac{\partial \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} = \frac{\partial \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}} + \frac{\partial \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \hat{\alpha}_i(\boldsymbol{\theta})} \cdot \frac{\partial \hat{\alpha}_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ ,  $d_{\theta\alpha_i} = \frac{\partial^2 \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta} \partial \alpha_i}$ . Making use of McCullah's (1987) asymptotic expansion for

$(\hat{\alpha}_i(\boldsymbol{\theta}) - \alpha_i)$  and  $(\hat{\alpha}_i(\boldsymbol{\theta}) - \alpha_i)^2$ , Carro (2007) derives the bias-corrected estimator from the modified score function of  $d_{\theta_i} = \frac{\partial \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \boldsymbol{\theta}}$ ,

$$\begin{aligned} \sum_{i=1}^N d_{\theta_i}^* &= \sum_{i=1}^N \left\{ d_{\theta_i}(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})) - \frac{1}{2} \frac{1}{d_{\alpha_i \alpha_i}(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))} \left( d_{\theta \alpha_i}(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})) \right. \right. \\ &\quad \left. \left. + d_{\alpha_i \alpha_i}(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta})) \frac{\partial \hat{\alpha}_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right) \right. \\ &\quad \left. + \frac{\partial}{\partial \alpha_i} \left( \frac{1}{E[d_{\alpha_i \alpha_i}(\boldsymbol{\theta}, \alpha_i)]} E[d_{\theta \alpha_i}(\boldsymbol{\theta}, \alpha_i)] \right) \Big|_{\hat{\alpha}_i(\boldsymbol{\theta})} \right\} = \mathbf{0}, \end{aligned} \quad (6.6.16)$$

where  $d_{\alpha_i \alpha_i} = \frac{\partial^2 \ell_i(\boldsymbol{\theta}, \alpha_i(\boldsymbol{\theta}))}{\partial \alpha_i^2}$  and  $d_{\alpha_i \alpha_i} = \frac{\partial^3 \ell_i(\boldsymbol{\theta}, \hat{\alpha}_i(\boldsymbol{\theta}))}{\partial \hat{\alpha}_i^3(\boldsymbol{\theta})}$ . Carro (2007) shows that the bias of the modified MLE,  $\hat{\boldsymbol{\theta}}^*$ , is also of order  $(\frac{1}{T^2})$  and has the same asymptotic variance as the MLE. His Monte Carlo studies show that the bias of the modified MLE is small with  $T = 8$ .

### 6.6.2 Bounding Parameters

When  $y_{i0}$  and  $\alpha_i$  are treated as random, the joint likelihood of  $f(\mathbf{y}_i, y_{i0} | \mathbf{x}_i)$  can be written in the form of conditional density of  $f(\mathbf{y}_i | y_{i0}, \mathbf{x}_i)$  times the marginal density  $f(y_{i0} | \mathbf{x}_i)$ ,

$$\begin{aligned} f(\mathbf{y}_i, y_{i0} | \mathbf{x}_i) &= \int f(\mathbf{y}_i | y_{i0}, \mathbf{x}_i, \alpha_i) f(y_{i0} | \mathbf{x}_i, \alpha_i) G(\alpha_i | \mathbf{x}_i) d\alpha_i, \\ i &= 1, \dots, N, \end{aligned} \quad (6.6.17)$$

where  $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$ ,  $\mathbf{x}_i = (\mathbf{x}'_{i1}, \dots, \mathbf{x}'_{iT})$ , and  $G(\alpha_i | \mathbf{x}_i)$  denote the conditional density of  $\alpha_i$  given  $\mathbf{x}_i$ . For model (6.6.8) with  $u_{it}$  following a standard normal distribution,  $N(0, 1)$ ,

$$\begin{aligned} f(\mathbf{y}_i | y_{i0}, \mathbf{x}_i, \alpha_i) &= \prod_{t=1}^T [\Phi(\mathbf{x}'_{it} \boldsymbol{\beta} + y_{i,t-1} \gamma + \alpha_i)]^{y_{it}} \\ &\quad \times [1 - \Phi(\mathbf{x}'_{it} \boldsymbol{\beta} + y_{i,t-1} \gamma + \alpha_i)]^{1-y_{it}} \\ i &= 1, \dots, N. \end{aligned} \quad (6.6.18)$$

If  $u_{it}$  follows a logistic distribution

$$\begin{aligned} f(\mathbf{y}_i | y_{i0}, \mathbf{x}_i, \alpha_i) &= \prod_{t=1}^T \frac{\exp[(\mathbf{x}'_{it} \boldsymbol{\beta} + y_{i,t-1} \gamma + \alpha_i)]^{y_{it}}}{1 + \exp(\mathbf{x}'_{it} \boldsymbol{\beta} + y_{i,t-1} \gamma + \alpha_i)}, \\ i &= 1, \dots, N. \end{aligned} \quad (6.6.19)$$

When  $G(\alpha | \mathbf{x})$  and the initial distribution  $f(y_0 | \mathbf{x})$  are known,  $\boldsymbol{\theta} = (\boldsymbol{\beta}, \gamma)$  can be estimated by the MLE. However,  $G(\alpha | \mathbf{x})$  and  $f(y_0 | \mathbf{x})$  are usually unknown. Although in principle, one can still maximize (6.6.17), the usual regularity conditions for the consistency of the MLE (e.g. Kiefer and Wolfowitz 1956) are violated because  $G(\alpha | \mathbf{x})$  is infinite dimensional (Cosslett 1981).

If  $\mathbf{x}$  and  $\alpha$  are discrete, Honoré and Tamer (2006) suggest a linear program approach to provide the bound of  $\boldsymbol{\theta}$ . Let  $A_j = (d_{j1}, \dots, d_{jT})$  be the  $1 \times T$  sequence of binary variables  $d_{jt}$ . Let  $\mathcal{A}$  denote the set of all  $2^T$  possible sequence of zeros and ones,  $A_j$ . Let

$P(y_{i0} | \mathbf{x}_i, \alpha_i)$  denote the probability of  $y_{i0} = 1$  given  $\mathbf{x}_i$  and  $\alpha_i$  and let  $f_0(\alpha, \mathbf{x})$  denote the distribution of  $y_{i0}$  given  $\mathbf{x}$  and  $\alpha$ . Then, conditional on  $P(y_{i0} | \mathbf{x}_i, \alpha_i)$ ,

$$\begin{aligned} f(\mathbf{y}_i | \mathbf{x}_i, f_0(y_{i0} | \mathbf{x}_i, \alpha_i), \alpha_i) &= P(y_{i0} | \mathbf{x}_i, \alpha_i) f(\mathbf{y}_i | y_{i0} = 1, \mathbf{x}_i, \alpha_i) \\ &\quad + (1 - P(y_{i0} | \mathbf{x}_i, \alpha_i)) f(\mathbf{y}_i | y_{i0} = 0, \mathbf{x}_i, \alpha_i), \end{aligned} \quad (6.6.20)$$

and

$$f(\mathbf{y}_i | \mathbf{x}_i, f_0(\cdot, \cdot), \boldsymbol{\theta}) = \int f(\mathbf{y}_i | \mathbf{x}_i, \alpha, f_0(\cdot, \cdot), \boldsymbol{\theta}) dG(\alpha | \mathbf{x}_i). \quad (6.6.21)$$

Let  $\pi(A | \mathbf{x}, \alpha; f_0(\cdot, \cdot), \boldsymbol{\theta})$  and  $P(A | \mathbf{x})$  be the probability of an event  $A$  in  $\mathcal{A}$  given  $(\mathbf{x}, \alpha)$  predicted by the model and the probability of an event  $A$  occurs given  $\mathbf{x}$ , respectively. Then  $\pi(A | \mathbf{x}, f_0(\cdot, \cdot), \boldsymbol{\theta}) = \int \pi(A | \mathbf{x}, \alpha; f_0(\cdot, \cdot), \boldsymbol{\theta}) dG(\alpha | \mathbf{x})$ . Define the set of  $(f_0(\cdot, \cdot), \boldsymbol{\theta})$  that is consistent with a particular data-generating process with probabilities  $\mathcal{P}(\mathcal{A} | \mathbf{x})$  as

$$\Psi = \left\{ (f_0(\cdot, \cdot), \boldsymbol{\theta}) : P[\pi(\mathcal{A} | \mathbf{x}; f_0(\cdot, \cdot), \boldsymbol{\theta}) = P(\mathcal{A} | \mathbf{x})] = 1 \right\}. \quad (6.6.22)$$

Then the bound of  $\boldsymbol{\theta}$  is given by

$$\Theta = \left\{ \boldsymbol{\theta} : \exists f_0(\cdot, \cdot) \text{ such that } P[\pi(\mathcal{A} | \mathbf{x}; f_0(\cdot, \cdot), \boldsymbol{\theta}) = P(\mathcal{A} | \mathbf{x})] = 1 \right\}. \quad (6.6.23)$$

Suppose  $\alpha$  has a discrete distribution with a known maximum number of points of support,  $M$ . The points of support are denoted by  $a_m$  and the probability of  $\alpha_i = a_m$  given  $\mathbf{x}$  is denoted by  $\rho_{mx}$ . Then

$$\begin{aligned} \pi(\mathcal{A} | f_0(\cdot, \cdot), \mathbf{x}, \boldsymbol{\theta}) &= \sum_{m=1}^M \rho_{mx} \left[ f_0(a_m, \mathbf{x}) \pi(\mathcal{A} | y_0 = 1, \boldsymbol{\theta}, \mathbf{x}; a_m) \right. \\ &\quad \left. + (1 - f_0(a_m, \mathbf{x})) \pi(\mathcal{A} | y_0 = 0, \boldsymbol{\theta}, \mathbf{x}; a_m) \right] \\ &= \sum_{m=1}^M z_{mx} \pi(\mathcal{A} | y_0 = 1, \mathbf{x}, \boldsymbol{\theta}; a_m) \\ &\quad + \sum_{m=1}^M z_{M+m, x} \pi(\mathcal{A} | y_0 = 0, \mathbf{x}, \boldsymbol{\theta}; a_m), \end{aligned} \quad (6.6.24)$$

where  $z_{mx} = \rho_{mx} f_0(a_m, \mathbf{x})$  and  $z_{M+m, x} = \rho_{mx} [1 - f_0(a_m, \mathbf{x})]$  for  $m = 1, \dots, M$ . The identified set  $\Theta$ , consists of the value of  $\boldsymbol{\theta}$  for which the following equations have a solution for  $\{z_{mx}\}_{m=1}^{2M}$ :

$$\begin{aligned} \sum_{m=1}^M z_{mx} \pi(A | y_0 = 1, \mathbf{x}, \boldsymbol{\theta}; a_m) + \sum_{m=1}^M z_{M+m, x} \pi(A | y_0 = 0, \mathbf{x}, \boldsymbol{\theta}; a_m) \\ = P(A | \mathbf{x}), \end{aligned} \quad (6.6.25)$$

and for all  $A \in \mathcal{A}$ ,

$$\sum_{m=1}^{2M} z_{mx} = 1, z_{mx} \geq 0. \quad (6.6.26)$$

Equations (6.6.25) and (6.6.26) have exactly the same structure as the constraints in a linear programming problem, so checking whether a particular  $\boldsymbol{\theta}$  belongs to  $\Theta$  can be

done in the same way that checks for a feasible solution in a linear programming problem provided  $P(A \mid \mathbf{x})$  can be consistently estimated. Therefore, Honoré and Tamer (2006) suggest bounding  $\theta$  by considering the linear programming problem:

$$\max_{\{z_{mx}, \{v_{jx}\}\}} \sum_j -v_{jx} \quad (6.6.27)$$

where

$$\begin{aligned} v_{jx} = & P(A_j \mid \mathbf{x}) - \sum_{m=1}^M z_{mx} \pi(A_j \mid y_{i0} = 1, \mathbf{x}, \theta; a_m) \\ & - \sum_{m=1}^M z_{M+m, x} \pi(A_j \mid y_{i0} = 0, \mathbf{x}, \theta; a_m) \end{aligned}$$

$$\text{for all } A_j \in \mathcal{A}, j = 1, \dots, 2^T, \quad (6.6.28)$$

$$1 - \sum_{m=1}^{2M} z_{mx} = v_{0x}, \quad (6.6.29)$$

$$z_{mx} \geq 0, \quad (6.6.30)$$

$$v_{jx} \geq 0. \quad (6.6.31)$$

The optimal function value for (6.6.27) is zero if and only if all  $v_{jx} = 0$ , i.e., if a solution exists to (6.6.25) and (6.6.26). If (6.6.25) and (6.6.26) do not have a solution, the maximum function value in (6.6.27) is negative. Following Manski and Tamer (2002), it can be shown that a consistent estimator of the identified region can be constructed by checking whether, for a given  $\theta$ , the sample objective function is within  $\epsilon$  of the maximum value of zero where  $P(A)$  is substituted by its consistent estimator. Since  $\mathbf{x}$  is discrete, one can mimic this argument for each value in the support of  $\mathbf{x}_i$ , which will then contribute a set of constraints to the linear programming problem.

### 6.6.3 Approximate Model

The dynamic logit model (6.5.26) implies that the conditional distribution of a sequence of response variables,  $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})'$  given  $\alpha_i, \mathbf{x}_i$ , and  $y_{i0}$ , can be expressed as

$$P(\mathbf{y}_i \mid \mathbf{x}_i, \alpha_i, y_{i0}) = \frac{\exp(\alpha_i \sum_{t=1}^T y_{it} + \sum_{t=1}^T y_{it}(\mathbf{x}'_{it}\boldsymbol{\beta}) + y_{i*}\gamma)}{\sum_{t=1}^T [1 + \exp(\mathbf{x}'_{it}\boldsymbol{\beta} + y_{i,t-1}\gamma + \alpha_i)]}, i = 1, \dots, N, \quad (6.6.32)$$

where  $y_{i*} = \sum_{t=1}^T y_{i,t-1}y_{it}$ .

The Honoré and Kyriazidou (2000) conditional approach discussed in Section 6.5 requires a specification of a suitable kernel function and the bandwidth parameters to weigh the response configuration of each subject in the sample on the basis of exogenous explanatory variables in which only exogenous variables are close to each other receiving large weights, implying a substantial reduction of the rate of convergence of the estimator to the true parameter value. Moreover, conditional on the certain configurations that leads

to a response function in terms of time changes of the covariates, implies the exclusion of time-invariant variables. Noting that the dynamic logit model (6.6.32) implies that the conditional log-odds ratio between  $(y_{it}, y_{i,t-1})$  is equal to

$$\log \frac{P(y_{it} = 0 \mid \alpha_i, \mathbf{x}_i, y_{i,t-1} = 0) \cdot P(y_{it} = 1 \mid \alpha_i, \mathbf{x}_i, y_{i,t-1} = 1)}{P(y_{it} = 0 \mid \alpha_i, \mathbf{x}_i, y_{i,t-1} = 1) \cdot P(y_{it} = 1 \mid \alpha_i, \mathbf{x}_i, y_{i,t-1} = 0)} = \gamma, \quad (6.6.33)$$

Bartolucci and Nigro (2010) suggest using the Cox (1972) quadratic exponential model to approximate (6.6.32)<sup>23</sup>

$$\begin{aligned} P^*(\mathbf{y}_i \mid \mathbf{x}_i, y_{i0}, \delta_i) \\ = \frac{\exp [\delta_i (\sum_{t=1}^T y_{it}) + \sum_{t=1}^T y_{it} (\mathbf{x}'_{it} \boldsymbol{\phi}_1) + y_{iT} (\psi + \mathbf{x}'_{iT} \boldsymbol{\phi}_2) + y_{iT} \tau]}{\sum_{d_{ij}} \exp [\delta_i (\sum_{t=1}^T d_{ijt}) + \sum_{t=1}^T d_{ijt} (\mathbf{x}'_{it} \boldsymbol{\phi}_1) + d_{iT} (\psi + \mathbf{x}'_{iT} \boldsymbol{\phi}_2) + d_{iT}^* \tau]} \end{aligned} \quad (6.6.34)$$

where  $\mathbf{d}_{ij} = (d_{ij1}, \dots, d_{ijT})$  denote the  $j$ th possible binary response sequence, and  $\sum_{d_{ij}}$  denotes the sum over all possible response sequence of  $\mathbf{d}_{ij}$ , such that  $\sum_{t=1}^T d_{ijt} = \sum_{t=1}^T y_{it}$ , and  $d_{ij}^* = d_{ij1} y_{i0} + \sum_{t=1}^T d_{ijt} d_{ij,t-1}$ . Model (6.6.34) implies that

$$P^*(y_{it} \mid \mathbf{x}_i, \delta_i, y_{i0}, \dots, y_{i,t-1}) = \frac{\exp \{y_{it} [\delta_i + \mathbf{x}'_{it} \boldsymbol{\phi}_1 + y_{i,t-1} \tau + e_t^*(\delta_i, \mathbf{x}_i)]\}}{1 + \exp [\delta_i + \mathbf{x}'_{it} \boldsymbol{\phi}_1 + y_{i,t-1} \tau + e_t^*(\delta_i, \mathbf{x}_i)]}, \quad (6.6.35)$$

where, for  $t < T$ ,

$$\begin{aligned} e_t^*(\delta_i, \mathbf{x}_i) &= \log \frac{1 + \exp [\delta_i + \mathbf{x}'_{i,t+1} \boldsymbol{\phi}_1 + e_{t+1}^*(\delta_i, \mathbf{x}_i) + \tau]}{1 + \exp [\delta_i + \mathbf{x}'_{i,t+1} \boldsymbol{\phi}_1 + e_{t+1}^*(\delta_i, \mathbf{x}_i)]} \\ &= \log \frac{P(y_{i,t+1} = 0 \mid \delta_i, \mathbf{x}_i, y_{it} = 0)}{P(y_{i,t+1} = 0 \mid \delta_i, \mathbf{x}_i, y_{it} = 1)}. \end{aligned} \quad (6.6.36)$$

The correction term (6.6.36) depends on future covariates. For the last period, it is approximated by

$$e_T^*(\delta_i, \mathbf{x}_i) = \psi + \mathbf{x}'_{iT} \boldsymbol{\phi}_2. \quad (6.6.37)$$

Model (6.6.35) may be viewed as a latent response model of the form

$$y_{it}^* = \mathbf{x}'_{it} \boldsymbol{\phi}_1 + \delta_i + y_{i,t-1} \tau + e_t^*(\delta_i, \mathbf{x}_i) + \eta_{it}, \quad (6.6.38)$$

with the logistically distributed stochastic term  $\eta_{it}$ . The correction term  $e_t^*(\delta_i, \mathbf{x}_i)$  may be interpreted as a measure of the effect of the present choice  $y_{it}$  on the expected utility (or propensity) at period  $(t+1)$ . The parameter  $\tau$  for the state dependence is the log-odds ratio between any pairs of variables  $(y_{i,t-1}, y_{it})$ , conditional on all the other response variables or marginal with respect to these variables.

The difference between the approximate model (6.6.34) and the dynamic logit model (6.6.32) is in the denominator. The former does not depend on the actual sequence  $\mathbf{y}_i$ , while the latter does. The advantage of model (6.6.34) or (6.6.35) is that the parameters for the unobserved heterogeneity,  $\delta_i$ , can be eliminated by conditioning on the sum of response variables over time just like the static logit model (6.3.15). When  $y_{i0}$  are observable, the

<sup>23</sup> To differentiate the approximate model from the dynamic logit model (6.6.32), we use  $\delta_i$  to represent the individual-specific effects and  $\boldsymbol{\phi}_1$  to represent the coefficients of  $\mathbf{x}_{it}$  in the approximate model (6.6.34).

structural parameters can be estimated by the conditional maximum likelihood estimator, as discussed in (6.3.22) when  $T \geq 2$ .

The relations between (6.6.32) and (6.6.34) can be seen through a Taylor series expansion of the nonlinear term of the logarithm of the dynamic logit model (6.6.32) at  $\alpha_i = \tilde{\alpha}_i, \beta = \tilde{\beta}$  and  $\gamma = 0$ ,

$$\begin{aligned} \sum_{t=1}^T \log [1 + \exp (\mathbf{x}'_{it} \beta + y_{i,t-1} \gamma + \alpha_i)] \\ \simeq \sum_{t=1}^T \{ \log [1 + \exp (\mathbf{x}'_{it} \tilde{\beta} + \tilde{\alpha}_i)] + \tilde{q}_{it} \\ \cdot [\mathbf{x}'_{it} (\beta - \tilde{\beta}) + (\alpha_i - \tilde{\alpha}_i)] \} + \tilde{q}_{i1} y_{i0} \gamma + \sum_{t=1}^T \tilde{q}_{it} y_{i,t-1} \gamma, \end{aligned} \quad (6.6.39)$$

where

$$\tilde{q}_{it} = \frac{\exp (\tilde{\alpha}_i + \mathbf{x}'_{it} \tilde{\beta})}{1 + \exp (\tilde{\alpha}_i + \mathbf{x}'_{it} \tilde{\beta})}. \quad (6.6.40)$$

Substituting (6.6.39) into the logarithm of (6.6.32) and renormalizing the exponential of the resulting expression yields the approximate model as

$$\begin{aligned} P^*(\mathbf{y}_i | \mathbf{x}_i, \alpha_i, y_{i0}) \\ = \frac{\exp (\alpha_i (\sum_{t=1}^T y_{it}) + \sum_{t=1}^T y_{it} (\mathbf{x}'_{it} \beta) - \sum_{t=2}^T \tilde{q}_{it} y_{i,t-1} \gamma + y_{i*} \gamma)}{\sum_{d_i} \exp [\alpha_i (\sum_{t=1}^T d_{ijt}) + \sum_{t=1}^T d_{ijt} (\mathbf{x}'_{it} \beta) - (\sum_{t=2}^T \tilde{q}_{it} d_{i,t-1}) \gamma + d_{ij*} \gamma]}, \\ i = 1, \dots, N. \end{aligned} \quad (6.6.41)$$

When  $\gamma$  is indeed equal to zero, the true model and the approximating model coincide. Both become the static logit model (6.3.14). The approximating model (6.6.34) or (6.6.41) implies that the conditional logit of  $y_{it}$  given  $\mathbf{x}'_i, \alpha_i$  and  $y_{i0}, \dots, y_{i,t-1}$ , is equal to

$$\begin{aligned} \log \frac{P^*(y_{it} = 1 | \mathbf{x}_i, \alpha_i, y_{i0}, \dots, y_{i,t-1})}{P^*(y_{it} = 0 | \mathbf{x}_i, \alpha_i, y_{i0}, \dots, y_{i,t-1})} \\ = \begin{cases} \mathbf{x}'_{it} \beta + y_{i,t-1} \gamma + \alpha_i + e_t(\alpha_i, \mathbf{x}_i) - \tilde{q}_{i,t+1} \gamma, & \text{if } t < T, \\ \mathbf{x}'_{it} \beta + y_{i,t-1} \gamma + \alpha_i, & \text{if } t = T, \end{cases} \end{aligned} \quad (6.6.42)$$

where

$$\begin{aligned} e_t(\alpha_i, \mathbf{x}_i) &= \log \frac{P^*(y_{i,t+1} = 0 | \mathbf{x}_i, \alpha_i, y_{it} = 0)}{P^*(y_{i,t+1} = 0 | \mathbf{x}_i, \alpha_i, y_{it} = 1)} \\ &= \tilde{q}_{i,t+1} \gamma. \end{aligned} \quad (6.6.43)$$

Equation (6.6.42) implies that

$$\log \frac{P^*(y_{it} = 1 | \mathbf{x}_i, \alpha_i, y_{i,t-1} = 1)}{P^*(y_{it} = 0 | \mathbf{x}_i, \alpha_i, y_{i,t-1} = 1)} - \log \frac{P^*(y_{it} = 1 | \mathbf{x}_i, \alpha_i, y_{i,t-1} = 0)}{P^*(y_{it} = 0 | \mathbf{x}_i, \alpha_i, y_{i,t-1} = 0)} = \gamma. \quad (6.6.44)$$



$$P^* \left( \mathbf{y}_i \mid \mathbf{x}_i, \alpha_i, y_{i0}, \sum_{t=1}^T y_{it} \right) = \frac{\exp \left\{ \sum_{t=2}^T y_{it} (\mathbf{x}_{it} - \mathbf{x}_{i,t-1}) \boldsymbol{\beta} - \sum_{t=2}^T \tilde{q}_{it} y_{i,t-1} \gamma + y_{i*} \gamma \right\}}{\sum_{d_i} \exp \left\{ \sum_{t=2}^T d_{ijt} (\mathbf{x}_{it} - \mathbf{x}_{i,t-1})' \boldsymbol{\beta} - \sum_{t=2}^T \tilde{q}_{it} d_{ij,t-1} \gamma + d_{ij*} \gamma \right\}},$$

$$i = 1, \dots, N. \quad (6.6.45)$$

To obtain the pseudo conditional MLE of the pseudo likelihood function (6.6.45), Bartolucci and Nigro (2010) suggest first assuming there was no state dependence ( $\gamma = 0$ ) and maximizing the conditional log-likelihood of the static logit model (6.3.22) for those  $i$  where  $0 < \sum_{t=1}^T y_{it} < T$  to obtain a preliminary estimate  $\tilde{\beta}$ . Then substituting  $\tilde{\beta}$  into (6.6.44) to obtain the revised pseudo-conditional MLE of  $\beta$  and  $\gamma$  through the Newton–Raphson iterative procedure. Their Monte Carlo studies show that the pseudo-conditional MLE has a very low bias for data generated by a dynamic logit model.