

Dynamic binary outcome models with maximal heterogeneity*

Martin Browning

Department of Economics

University of Oxford

Martin.Browning@economics.ox.ac.uk

Jesus M. Carro

Departamento de Economía,

Universidad Carlos III de Madrid.

jcarro@eco.uc3m.es

First Draft: July 2007[†]

This Draft: January 2011

JEL classification: C23, C24, J64

Keywords: discrete choice, Markov processes, nonparametric identification, unemployment dynamics.

Abstract

Most econometric schemes to allow for heterogeneity in micro behaviour have two drawbacks: they do not fit the data and they rule out interesting economic models. In this paper we consider the time homogeneous first order Markov (HFOM) model that allows for maximal heterogeneity. That is, the modelling of the heterogeneity does not impose anything on the data (except the HFOM assumption for each agent) and it allows for any theory model (that

*For comments and useful suggestions, we thank Enrique Sentana, Whitney Newey, Ivan Fernandez-Val, Sara Ayo, and participants at seminars at Boston University; MIT/Harvard; Yale University; Nuffield (Oxford); IFS (London); CEMFI; Manchester; Columbia, CAM (Copenhagen) and a conference at the Tinbergen Institute. The second author gratefully acknowledges that this research was supported by a Marie Curie International Outgoing Fellowship within the 7th European Community Framework Programme.

[†]The first draft was titled "Identification of the dynamic discrete choice model" and was presented at the CAM Summer Workshop (University of Copenhagen) in July 2007.

gives a HFOM process for an individual observable variable). ‘Maximal’ means that the joint distribution of initial values and the transition probabilities is unrestricted.

We establish necessary and sufficient conditions for generic local point identification of our heterogeneity structure and show how it depends on the length of the panel. A feasible ML estimation procedure is developed. Tests for a variety of subsidiary hypotheses such as the assumption that marginal dynamic effects are homogeneous are developed.

We apply our techniques to a long panel of Danish workers who are very homogeneous in terms of observables. We show that individual unemployment dynamics are very heterogeneous, even for such a homogeneous group. We also show that the impact of cyclical variables on individual unemployment probabilities differs widely across workers. Some workers have unemployment dynamics that are independent of the cycle whereas others are highly sensitive to macro shocks.

1. Introduction.

Models with a binary outcome that depends in part on previous realizations of the outcome - dynamic binary outcome models - are common in applied microeconomics. Some examples include: labour force participation (Heckman (1981), Hyslop (1999)); smoking (Becker *et al* (1994)); firms exporting (Bernard and Jensen (2004)); stock market participation (Alessie *et al* (2004)) and taking up a welfare program (Gottschalk and Moffitt (1994) and Ham and Shore-Sheppard (2005)). The usual time-homogeneous first order Markov model for unit i ($= 1, ..N$) in period t ($t = 0, ..T$) is:

$$\Pr(y_{it} = 1 \mid y_{i,t-1}, x_{it}) = F(\eta_i + \alpha y_{i,t-1} + \beta x_{it}) \quad (1.1)$$

where $F(\cdot)$ is a probability distribution function and y_{it} is a binary variable indicating, for example, that person i had some unemployment in period t . This ‘linear index model’ which only allows for a heterogeneous ‘intercept’ η_i is widely used but it does have problems; Browning and Carro (2007) discuss these but it is worth repeating the objections.

The first problem is that the imposition of common slope parameters (α and β) restricts the class of structural models that are consistent with the reduced form (1.1). For example, consider two people, a and b , with the same value of the x variables (so we can ignore them), and for whom a has a lower probability of being unemployed if

they were employed in the previous year:

$$F(\eta_a) < F(\eta_b) \quad (1.2)$$

For example, a might choose a ‘safer’ job than b . Now suppose we impose the ‘same slope’ homogeneity assumption $\alpha_a = \alpha_b = \alpha$. This implies:

$$F(\eta_a + \alpha) < F(\eta_b + \alpha) \quad (1.3)$$

This rules out, for example, that a ’s caution leads her to spend more time looking for a ‘safe’ job, so that her probability of remaining unemployed is *higher* than b ’s. Thus the choice of a statistical scheme for dealing with heterogeneity has substantive restrictions on the set of admissible structural models.

The second problem with the conventional approach is that whenever we have long enough panels to estimate the model for each unit individually with minimal bias, we do find substantial heterogeneity in the both the ‘intercept’ and ‘slope’ parameters in (1.1). A situation where this is the case can be found in Browning and Carro (2010). Additional evidence will be provided in the empirical illustration in this paper.

Model (1.1) with maximal heterogeneity has:

$$\Pr(y_{it} = 1 \mid y_{i,t-1}, x_{it}) = F(\eta_i + \alpha_i y_{i,t-1} + \beta_i x_{it}) \quad (1.4)$$

In addition to the homogeneity restrictions, model (1.1) is imposing two kind of parametric restrictions: the parametric form implied by the linear index and the probability distribution function $F(\cdot)$. In this paper, we consider not only a semiparametric form but also the nonparametric case as well as having maximal heterogeneity all throughout the paper.¹ A nonparametric time-homogeneous first order Markov process with maximal heterogeneity will look directly at the transition probabilities allowing them to be different for each individual:

$$\Pr(y_{it} = 1 \mid y_{i,t-1}, x_{it}) = p_{ixy_{t-1}} \quad (1.5)$$

where we have one parameter to be estimated for each i and value of x and the lag of y . This does not impose any restrictions on the structural model (except, of course, for the assumption of time invariance and no effects higher than the first order that define the model considered in this paper) and it will fit any data that is generated by a time-homogeneous first order Markov process (HFOM). For the simpler case

¹Notice also that in (1.1) an extra homogeneity assumption is imposed by assuming all i have the same $F(\cdot)$. In our nonparametric approach this homogeneity assumption is not imposed either.

without x variables there is a one to one correspondence between (1.4) and (1.5) and, therefore, any $F(\cdot)$ will give the same transition probabilities. For the general case with x variables, a semiparametric form assuming a function $F(\cdot)$ in (1.4) will impose some parametric restrictions that are not imposed in (1.5).

Identifying and estimating the whole set of transition probabilities in (1.5) - the whole set of parameters if we consider (1.4) - or their distribution over the population, allows us to obtain any parameter of interest in this problem, including the average marginal effect (also known as average partial effect, APE) of a explanatory variable over the outcome y_{it} . This is important since different studies and questions require us to obtain different parameters of interest. Moreover, the average may not be a very informative measure because of the discrete nature of the problem. For instance, the APE could be found to be very small only because of a group in the population for which a change of a variable does not have enough effect as to change their y_{it} given their other observable and unobservable circumstances. In this case the APE will not be informative about other parts of the population for which the impact can be very large because they are close to the margin that make them change their y_{it} . In this situation measures such as the median marginal effect are more informative. Also, even if we look only at mean effects, there is more than one that could be of interest: the mean effect for a randomly drawn individual (see Chamberlain, 1984) or ATE in the treatment effect literature, the average marginal effect of x when $x = x_1$ only for those with $x = x_1$ (see Altonji and Matzkin, 2005), the ‘average treatment on treated’, etc.. Furthermore, identifying and estimating the whole HFOM model will allow to obtain the entire distribution in the population of the effect of a variable over the outcome. In a program evaluation context, Heckman, Smith and Clements (1997) present situations in which the entire distribution, and not only the mean effect, is the policy parameter of interest. In the IO literature it is also of interest to identify the entire distribution of the individual price elasticities when estimating demand functions; see for example Nevo (2001).

Given the difficulties in estimating (1.1) with small and fixed T (see Arellano and Honoré (2001)), tackling (1.5) or (1.4) is a formidable task. In Browning and Carro (2010) we suggested two estimation methods for the simple case without x variables, that rely on reducing the bias or RMSE for estimates based on each unit. This gives estimates for each unit and then the distribution for (η, α) can be taken as the empirical distribution of these estimates (or some smoothed version of it).

In Browning and Carro (2010), identification and estimation of (1.5) without imposing any restriction on the distribution of (η, α) nor on the initial condition, relies on the T dimension; that is, it is only consistent when $T \rightarrow \infty$. In this paper we propose an alternative approach that relies on large N . In general the

model is not nonparametrically identified from a cross section of observations of fixed length T .² This negative result is our starting point in this paper: identification from the cross section is our goal since we typically do not have panels with a very large number of periods. Nevertheless, this negative result on identification does not imply that we cannot learn anything from a cross section of paths with a fixed T . In general, some restrictions will have to be imposed on the distribution of the heterogeneity to achieve point identification. The interesting question is the nature of the restrictions we have to impose, or how much information about our model with maximal heterogeneity we can identify from a cross section of length T . To answer this question we use finite discrete mixture distributions for the joint set of unknown heterogeneous parameters. We refer to this as the *nonparametric discrete scheme* since no restriction is imposed other than there is a finite and discrete number of points of support on this distribution.

An advantage of this discrete scheme is that it allows us to go from the full homogeneous case (one point of support) to the totally unrestricted case (as many points of support as N) within the same scheme. Also, given the discrete nature of problem and the finite number of possible observations, it is clear that we cannot nonparametrically identify a continuous distribution. So, the nonparametric discrete distribution is our route to study nonparametric identification. If the underlying distribution is thought to be continuous, then the discrete scheme can be seen as an approximation to the continuous distribution. Given the discrete nature of problem, any continuous distribution of the heterogeneity can be approximated with finite discrete mixture distributions. Alternatively, if we are willing to assume a parametric form for a continuous distribution, we show in section 3.1 how our identification analysis can be modified to cover this case.

The identification issue in this scheme will be: how many points of support can we take for a given T ? A major gain from looking at models identified from a cross section with fixed T is that there is no incidental parameters problem nor finite sample bias problem from not having a large number of periods.

Kasahara and Shimotsu (2009) take a different approach to a more general problem that includes the model we consider here, as well as other models. One of the examples included in their paper to illustrate their results is model (1.4) without x variables. However, for this case they do not give identification conditions for an arbitrary number of periods. For example, their most important result for this model (proposition 7 in Kasahara and Shimotsu, 2009) requires $T \geq 8$. Also they give stronger sufficient conditions than the conditions derived in this paper, whereas here

²In general, not even the restrictive model (1.1) with only one fixed effect is identified; see Honoré and Tamer (2006).

we derive sufficient and necessary conditions for identification. Moreover, their conditions are nontrivial to check in actual data, whereas our conditions are simple to check.

A different and interesting analysis is to look at set identification for the cases that are not point identified. In particular to derive bounds in the non-identified situation when no restriction or distribution is assumed for the heterogeneous parameters. Chernozhukov, Fernandez-Val, Hahn and Newey (2009) do this for the average marginal effect in models such as the ones considered here; they derive results showing that bounds can shrink and converge as T grows.

In sections 2 – 4 we study in detail the simpler dynamic HFOM model without x covariates. Studying the model without x covariates helps understanding the problem, and all the results derived for this case will be extended to the more interesting case with covariates that is taken up in section 5. Furthermore, the case without covariates will be a worst case reference in terms of identification; as we will show, having an exogenous x that is not constant across individuals facilitates identification. Sections 2 and 3 consider restrictions from the model and identification respectively. In section 4 we consider estimation and testing. In Section 6 we apply the techniques we develop to a long panel of Danish workers who are very homogeneous in terms of observables. Section 7 concludes.

The principal contributions of paper are:

- We provide necessary nonparametric conditions for any panel data set with binary outcomes to be consistent with a time-homogeneous first order Markov (HFOM) process. These conditions are simple and fast to check.
- Assuming the data has been generated by a HFOM process (both with and without covariates), we study identification for two types of distributions for the unobserved heterogeneity: parametric continuous and nonparametric discrete. In the latter case, it is shown that we can have a much richer distribution than the two point distribution usually found in applied work and still keep unrestricted important features of the distribution of the heterogeneity such as the initial condition or the correlation between the transition probabilities. Our main result provide necessary and sufficient conditions of local point identification in general.
- We give exact results on how identification depends on the length of the panel and on the covariates.
- We provide a framework that allows that macro variables have different effects for different agents.

2. HFOM model restrictions.

2.1. The research question.

We consider first a dynamic discrete choice model with no covariates in order to more easily study and understand the problem. The results derived for this case will be very useful for the case with covariates. The data consist of paths $\{y_{i0}, y_{i1}, \dots, y_{iT}\}_{i=1,2,\dots,N}$ where y_{it} is the value of a binary variable for unit i . We assume a time-homogeneous first order Markov (HFOM) process for each unit and define transition probabilities (1.5) in this case:

$$G_i = pr(y_{it} = 1 \mid y_{i,t-1} = 0) \quad (2.1)$$

$$H_i = pr(y_{it} = 1 \mid y_{i,t-1} = 1) \quad (2.2)$$

and the unconditional probability of a unit value for the initial observation:

$$P_i = pr(y_{i0} = 1) \quad (2.3)$$

This direct formulation is much more convenient to work with than the usual econometric specification given in (1.4) for two reasons. The first reason is that we do not have to specify any probability distribution function $F(\cdot)$, so we are nonparametric in modeling this HFOM. This reason does not have much consequences in this simpler model because allowing for maximal heterogeneity is enough to fit any data that is generated by a HFOM process when there is no x covariates. There is a one to one correspondence between (α_i, η_i) and (G_i, H_i) and, therefore, any F will give the same (G_i, H_i) transition probabilities. However in case with covariates the semiparametric form (1.4) will be imposing two kind of parametric restrictions: (i) the parametric form implied by the linear index and (ii) the probability distribution function $F(\cdot)$.

The second reason for this direct formulation is that parameters of (1.4) do not have any meaning on their own, apart from being different from zero or their sign. In contrast, (P_i, G_i, H_i) are probabilities and have a clear interpretation. Nevertheless the values of the parameters (P_i, G_i, H_i) are not usually of primary interest; rather they can be used to generate any other ‘outcomes or parameters of interest’. There are several candidates but the most widely considered for this model without covariates are the *marginal dynamic effects*:

$$\begin{aligned} M_i &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 1) - \Pr(y_{it} = 1 \mid y_{i,t-1} = 0) \\ &= H_i - G_i \end{aligned} \quad (2.4)$$

and the long run proportion of unit values:

$$\begin{aligned} L_i &= \frac{\Pr(y_{it} = 1 \mid y_{i,t-1} = 0)}{\Pr(y_{it} = 1 \mid y_{i,t-1} = 0) + \Pr(y_{it} = 0 \mid y_{i,t-1} = 1)} \\ &= \frac{G_i}{1 + G_i - H_i} \end{aligned} \quad (2.5)$$

Given that these values are heterogenous in i , their distribution over the population or some moments of them are the parameters of interest. An example, though not necessarily the most informative measure, is the *average marginal dynamic effect*:

$$E[M_i] = \int \int (H_i - G_i) dF_{(G,H)}(G_i, H_i) \quad (2.6)$$

where $F_{(G,H)}(G_i, H_i)$ is the joint distribution of G and H we want to identify. Another common object of interest is the probability that $y_{it} = 1$ in any given period t ; this is given by the Chapman-Kolmogorov equations applied to the initial probability and the transition probabilities. As explained in the introduction, there is more than one parameter of interest and identifying the whole HFOM model will allow to obtain any of them, including the entire distribution of M_i in the population.

Given this, our research question is: given a large- N , fixed- T panel, what can we (point) identify about the distribution of (P, G, H) over the population?

2.2. Enumerating paths.

For the moment we can drop the i subscript. There are $\Gamma = 2^{T+1}$ possible paths. The probability of a path j is given by:

$$p_j(P, G, H) = P^{y_0^j} (1 - P)^{(1-y_0^j)} G^{n_{01}^j} (1 - G)^{n_{00}^j} H^{n_{11}^j} (1 - H)^{n_{10}^j} \quad (2.7)$$

where n_{01}^j is the number of $0 \rightarrow 1$ transitions for path j and similarly for the other three transitions. We shall often use the $T = 2$ case to illustrate general points; Table 2.1 gives the probabilities for the eight possible paths. In all that follows we shall always order paths using a binary representation for ordering the elements for $t = 0, 2 \dots T$. Thus the first path is always 00..00, the second path is always 00..01 and the last path is always 11..11.

2.3. The general problem.

To consider the restrictions from the model and identification we assume that we are given population values for the probabilities of each of the Γ outcomes. Denote the population values by π_j for $j = 1, 2 \dots \Gamma$. Let (P, G, H) be distributed over $[0, 1]^3$

Case	Path	n_{00}	n_{01}	n_{10}	n_{11}	Probability of case j , p_j
1	000	2	0	0	0	$(1 - P)(1 - G)(1 - G)$
2	001	1	1	0	0	$(1 - P)(1 - G)G$
3	010	0	1	1	0	$(1 - P)G(1 - H)$
4	011	0	1	0	1	$(1 - P)GH$
5	100	1	0	1	0	$P(1 - H)(1 - G)$
6	101	0	1	1	0	$P(1 - H)G$
7	110	0	0	1	1	$PH(1 - H)$
8	111	0	0	0	2	PHH

Table 2.1: Outcomes for three periods (T=2)

with an unknown density $f(P, G, H)$. The population proportions are given by the integral equations:

$$\pi_j = \int_0^1 \int_0^1 \int_0^1 p_j(P, G, H) f(P, G, H) dPdGdH, \quad j = 1, 2, \dots, \Gamma \quad (2.8)$$

Since the p_j 's and the π_j 's sum to unity, $f(\cdot)$ will be a well defined density:

$$\begin{aligned} 1 &= \sum_{j=1}^{\Gamma} \pi_j = \int_0^1 \int_0^1 \int_0^1 \sum_{j=1}^{\Gamma} p_j(P, G, H) f(P, G, H) dPdGdH \\ &= \int_0^1 \int_0^1 \int_0^1 f(P, G, H) dPdGdH \end{aligned} \quad (2.9)$$

The econometric issues are:

1. Given a set of observed π_j 's for $j = 1, \dots, 2^{T+1}$, can we find a density function $f(P, G, H)$ such that (2.8) holds?
2. If we can find such a function for a given set of π_j 's, is it unique?
3. If we can find a unique inverse function, is the inverse mapping a continuous function of the values π_j ?

These are the usual set of conditions for a well posed inverse problem. The first condition asks if the model choice (in this case the form of the $p_j(P, G, H)$ functions due to the HFOM assumption) imposes any restrictions on observables. The second is the classical identification condition: given that the data are consistent with the model, can we recover unique estimates of the unknowns, in this case, the density $f(P, G, H)$. The final condition requires that the estimate of the unknown is 'stable' in the sense that small changes in the distribution of observables lead to small changes in the inferred unknowns. The continuity of the inverse mapping is also useful for

estimation since we can recover consistent estimates of the structural form (in this case, $f(\cdot)$) from consistent estimates of the reduced forms (the π_j 's).

2.4. Restrictions.

Turning to the first question, we ask whether any observed π_j 's that sum to unity could be generated by a HFOM process. The answer is clearly negative, since the data might have been generated by, for example, a time-homogeneous second order Markov scheme or a time-inhomogeneous first order process (or even more general models). Thus the time-homogeneity first order assumption will usually impose restrictions. The restrictions are a combination of equality restrictions and inequality restrictions. Considering (2.7) and (2.8) we have the following equality restrictions:

Lemma 2.1. *Given two paths j and j' , if*

$$y_0^j = y_0^{j'}, n_{00}^j = n_{00}^{j'}, n_{01}^j = n_{01}^{j'}, n_{10}^j = n_{10}^{j'}, n_{11}^j = n_{11}^{j'} \quad (2.10)$$

then $\pi_j = \pi_{j'}$.

Thus two population proportions will be equal if they have the initial value and the same number of transitions. For example, for $T = 3$ (that is, four periods of observation) the two paths 0010 and 0100 have the same initial value and the same number of transitions and hence the same probability,

$$\pi_{0010} = \pi_{0100} = \int_0^1 \int_0^1 \int_0^1 ((1 - P)(1 - G)HGf(P, G, H)) dPdGdH, \quad j = 1, 2, \dots, \Gamma \quad (2.11)$$

These are necessary conditions. There are further inequality restrictions. Consider, for example, the case of $T = 2$; see Table 2.1. There are no equality restrictions of the kind described in the Lemma. However, the restriction that $G \in [0, 1]$ imposes that

$$p_2(P, G, H) = (1 - P)(1 - G)G \leq 0.25 \quad (2.12)$$

Thus we have:

$$\pi_2 = \int_0^1 \int_0^1 \int_0^1 p_2(P, G, H) f(P, G, H) dPdGdH \leq 0.25 \quad (2.13)$$

Moreover, if π_2 is actually equal to 0.25 then $P = 0$ and $G = 0.5$ which in turn imposes $\pi_1 = 0.25$. Although we are not able to characterize the full set of necessary

# periods	T	$\Gamma = 2^{T+1}$	r_T	R_T
3	2	8	8	0
4	3	16	14	2
5	4	32	22	10
6	5	64	32	32
7	6	128	44	84
8	7	256	58	198
9	8	512	74	438
10	9	1024	92	932
11	10	2048	112	1936
12	11	4096	134	3962
13	12	8192	158	8034
14	13	16384	184	16200
15	14	32768	212	32556
16	15	65536	242	65294
24	23	$\sim 16.8 \times 10^6$	554	$\sim 16.8 \times 10^6$

Table 2.2: Numbers of possible paths, number of independent cases and number of restrictions

and sufficient conditions for a given π vector to be generated by a HFOM process, we show below how to test for them.

Using the Lemma above we can calculate the number of paths that are the same for any T , without considering the distribution $f(\cdot)$. For small T this calculation can be done by generating all the possible paths and counting with a computer. However, the following proposition gives a simple analytic formula for the number of different paths for any T , denoted by r_T .

Proposition 2.2. *The number of different paths in values of the vector $\pi = (\pi_1, \dots, \pi_j, \dots, \pi_T)'$ whose π_j elements are defined in (2.8) is*

$$r_T = T(T + 1) + 2 \tag{2.14}$$

The proof is given in Appendix A.1.

Table 2.2 presents the results for sample lengths of up to 16 and for 24 (the number used in our empirical example below). The values in the column headed r_T give the number of ‘independent’ values of the vector π and the column headed R_T gives the number of restrictions. For medium sized panels the reduction in the number of equations is quite dramatic. For example, for $T = 6$ we have 128 equations and 84 restrictions. This simply highlights that the first order and time-homogeneity assumptions impose strong restrictions if we have several periods of observations.

It is convenient to partition paths into groups based on their having the same probabilities. Define groups $k = 1, 2, \dots, r_T$ with $\pi_j = \pi_{j'}$ implying that j and j' are in the same group. Let n_k denote the number of members of group k and re-write (2.8) as:

$$\pi_k = n_k \int_0^1 \int_0^1 \int_0^1 p_k(P, G, H) f(P, G, H) dP dG dH, \quad k = 1, 2, \dots, r_T \quad (2.15)$$

Thus for $T = 5$, for example, we have 32 equations if the HFOM implications are not rejected. Below we shall present a maximum likelihood estimator for our model. When we do this, we shall show how to test for the restrictions implicit in the assumption that our finite sample data are generated by a HFOM process. We turn now to identification.

3. Identification.

Suppose the restrictions for the HFOM model developed in the previous section are not rejected. It is clear that with a finite set of path probabilities we cannot nonparametrically identify a continuous density $f(P, G, H)$ from the finite set of equations (2.15). If we had a continuous covariate and allowed that it had a homogeneous marginal effect on the parameters we could potentially identify the continuous distribution.³ Since we are here interested in identification without imposing arbitrary homogeneity schemes, this option is not open to us. This leaves us with two broad alternatives.

3.1. Identification of a parametric distribution.

The first broad alternative is take a known *parametric distribution function* $f(P, G, H; \beta)$ where β is an unknown L -vector. Thus:

$$\pi_k(\beta) = n_k \int_0^1 \int_0^1 \int_0^1 p_k(P, G, H) f(P, G, H; \beta) dP dG dH, \quad k = 1, 2, \dots, r_T \quad (3.1)$$

The identification issue is to ask whether we can identify the vector of parameters β . The Jacobian is the matrix:

$$J = \left[\frac{\partial \pi_k(\beta)}{\partial \beta_l} \right]_{k=1, \dots, r_T, l=1 \dots L} \quad (3.2)$$

³Subject to support restrictions that allow us to drive any probability to the limits of zero and unity.

In general we require that this matrix has a rank L , so that a necessary condition for (local) identification is $L \leq r_T$. For example, if we take a 9 parameter distribution for $f(P, G, H; \beta)$ then we could not point identify with $T = 2$ ($r_T = 8$) without imposing at least one restriction; for example that P is uncorrelated with (G, H) . If we take a mixture of two such distributions we have 19 parameters (the two sets of distributional parameters and the mixing probability) which would require $T \geq 4$. If we have a long panel then many components are allowed; for example, with $T = 23$ we could theoretically identify the parameters of a parametric model with 55 component nine parameter distributions. Given the order condition $L \leq r_T$, the rank of (3.2) would need to be checked for the particular parametric form chosen.

3.2. Identification for the nonparametric discrete scheme.

The second broad alternative assumption is that we have a *discrete finite mixture distribution* for (P, G, H) . For this, we consider nonparametric identification. We note that our use of a discrete distribution to capture heterogeneity is different to that suggested by Heckman and Singer (1984). They show that the distribution of a continuous latent variable is nonparametrically identified for a particular parametric duration model. They then suggest that the continuous distribution can be reasonably approximated by a discrete distribution with a small number of support points. In contrast, in our scheme the continuous distribution is *not* nonparametrically identified and the recourse to a discrete distribution is one route to nonparametric point identification. If, nonetheless, the underlying distribution is thought to be continuous, a discrete distribution can be taken as a good approximation to it because of the discrete nature of the problem; and our analysis would be used to know whether that nonparametric discrete distribution is identified or not. Of course there is always the possibility of considering a parametric continuous distribution; see section 3.1.

We take S distinct points of support $\{(P_1, G_1, H_1), \dots, (P_S, G_S, H_S)\}$ with probabilities given by the $(S \times 1)$ vector θ with non-negative individual values, θ_s , that sum to unity. The discrete analogue to (2.8) is:

$$\pi_j = \sum_{s=1}^S p_j(P_s, G_s, H_s) \theta_s \quad j = 1, 2, \dots, \Gamma \quad (3.3)$$

Define the $(\Gamma \times S)$ matrix A by:

$$A_{js} = p_j(P_s, G_s, H_s), \quad j = 1, 2, \dots, 2^{T+1}, \quad s = 1, 2, \dots, S \quad (3.4)$$

so that (2.15) can be written in matrix form as:

$$\pi = \mathbf{A}\theta \tag{3.5}$$

We take the support points and the probabilities to be unknown so that we have to solve for the values of $\{\mathbf{P}, \mathbf{G}, \mathbf{H}\}$ (the vectors of parameters) and θ . We refer to this as the *nonparametric discrete scheme*. An advantage of using this discrete distribution is that it allows us to go from the full homogeneous case (one point of support) to the totally unrestricted case (as many points of support as N) within the same scheme. Also, given the finite number of possible observations, any underlying continuous distribution of the heterogeneity can be approximated with a discrete finite mixture distribution. The identification issue is: how many periods we need to identify a distribution with S points of support?

Certainly not any discrete distribution with finite points of support will be identified from π . For example, it is easy to see that there are many distributions of $\{\mathbf{P}, \mathbf{G}, \mathbf{H}\}$ with $S = 8$ that will give the same proportions with $T = 2$.⁴ Therefore we cannot identify the distribution of $\{\mathbf{P}, \mathbf{G}, \mathbf{H}\}$ with $S = 8$, from the π we observe when $T = 2$. We need more periods to identify it.

From (3.5), for given S , we have a mapping from $(4S - 1)$ unobservables to observables given by:

$$\pi(\mathbf{P}, \mathbf{G}, \mathbf{H}, \theta_1, \dots, \theta_S) = \mathbf{A}(\mathbf{P}, \mathbf{G}, \mathbf{H})\theta$$

where the S -vector θ is normalized to sum to unity. In Appendix B.1 we show that identification in this system is equivalent to studying identification in the system conditional on the first observation. The reason is that from the distribution of $\{\mathbf{G}, \mathbf{H}\}$ conditional on y_0 plus the aggregate proportion of y_0 that we observe, $\Pr(y_0 = 1)$, there is a one to one map to the distribution of the heterogeneity in the distribution of $y_0, \{P_s\}_{s=1}^S$. This makes clear that we can only identify a finite discrete distribution of the heterogeneity on the initial observation as long as it is correlated with the distribution of $\{\mathbf{G}, \mathbf{H}\}$. If there were more heterogeneity on the initial observation but it were exogenous with respect to heterogeneity in the transition of the rest of the periods, we would not identify it since we only have one temporal realization of that independent distribution of y_0 . On the other hand we most often only care about the distribution of the initial condition as long as it is correlated with the distribution of the rest of the periods and ignoring it leads to misleading conclusions.

We denote the Jacobian matrix of the system conditional on the first observation

⁴To see this simply take two different sets of values of $\{\mathbf{P}, \mathbf{G}, \mathbf{H}\}$ with $S = 8$ such that \mathbf{A} is invertible and then there is a $\theta = \mathbf{A}^{-1}\pi$ for each set that defines two different sets of values of the parameters that imply the same π with $T = 2$.

S	1	2	3	4	5	6	7	8	14	25	50	100
$\min T + 1$	2	3	4*	5*	5*	6*	6*	6	8*	11*	15*	21*
r_T	4	8	14	22	22	32	32	32	58	112	212	422

*Over-identified

Table 3.1: Rank of the Jacobian and maximum number of points of support

by $\mathbf{J}_r(T, S)$. For local point identification we require that the rank of $\mathbf{J}_r(T, S)$ is greater than or equal to the number of parameters. In B we show that, in general (that is. except on a set of measure zero), the rank of \mathbf{J}_r is:

$$\text{rank}(\mathbf{J}_r) = \min(r_T - 2, 4S - 2) \quad (3.6)$$

This means that the parameters of S support points and their probabilities can only be point identified if the number of parameters is not greater than the rank of \mathbf{J}_r ; using (3.6), this requires:

$$S \leq \frac{r_T}{4} = \frac{T(T+1) + 2}{4} = \Upsilon_T \quad (3.7)$$

where the maximum S increases quadratically with T . From this condition we can calculate the minimum T needed to point identify a model with S number of points on support of the distribution of $\{\mathbf{P}, \mathbf{G}, \mathbf{H}\}$:

$$\min T = \left\lceil \frac{-1}{2} + \sqrt{\frac{-7}{4} + 4S} \right\rceil \quad (3.8)$$

where $\lceil x \rceil$ gives the smallest integer greater than or equal to x . Table (3.1) presents that minimum number of periods, $\min T + 1$, for some cases. As can be seen from Table (3.1), to identify a relatively rich distribution with 14 different points of support we only need a relatively short panel ($T = 7$). Even a short panel ($T = 4$, for example) is more than we need to identify a distribution with more than the two points commonly used in applied work.

Furthermore, this condition based on the rank of the Jacobian is not only sufficient for local identification but it is also necessary (almost everywhere) because they are regular points. The only non-regular points are those in the set of measure zero at which $\text{rank}(\mathbf{J}_r) < \min(r_T - 2, 4S - 2)$, because, as shown in appendix B.5, the rank is constant almost everywhere. All the previous results are summarized in the following proposition that gives sufficient and necessary conditions to locally identify our model almost everywhere:

Proposition 3.1. *The joint distribution of $\{\mathbf{P}, \mathbf{G}, \mathbf{H}\}$ with S points of support in*

system (3.5) is locally identified almost everywhere if and only if:

$$T \geq \frac{-1}{2} + \sqrt{\frac{-7}{4} + 4S} \quad (3.9)$$

The proof is given in Appendix B.

The condition in this proposition is weaker than the condition derived in Browning and Carro (2011), since we have shown that smaller T is required in general. Furthermore, the number of periods we need to identify the system in general increases at a root square rate with the number of points of support of the distribution we try to identify, as opposed to increasing our need of periods linearly (that is, at the same rate) as in Browning and Carro (2011). The reason for these differences is that in Browning and Carro (2011) we were not using moment conditions where both G and H interacted. The advantage of the approach in Browning and Carro (2011) is that the result there holds everywhere in the parameter space.

3.2.1. Non-regular Points

Given that the proof of Proposition 3.1 is based on the rank of the Jacobian, the set of points for which \mathbf{J}_r does not have full rank when $T = \min T$ is the set of particular points where Proposition 3.1 does not apply. We have shown that this set is of measure zero. However, we may want to know, firstly if there is any interesting case in that set. On the other hand, in this set of points, the condition on the rank of the Jacobian is sufficient but it is not necessary, since they are non-regular points. This means that it might be possible that they are identified even with T at which the rank of the Jacobian is smaller than the number of unknowns. Given this we would like to know also if this is the case or if a higher T is required. If they are identified even with T such that the rank of the Jacobian is smaller than the number of unknowns, then the condition (3.9) will be a sufficient condition also for non-regular points. If a higher T is required for identification, then condition (3.9) will be necessary but not sufficient. In all the non-regular points we have found, not only the Jacobian is not full rank (when considering $T = \min T$), but also the model is not identified; thus condition (3.9) about the number of periods is necessary but not sufficient in those non-regular points.

Given that there is no explicit solution to system (3.5) and it contains many non-linear equation and unknowns, we have worked with the simpler case $T = 2$, $S = 2$ to locate non-regular points and study its identification. Recall that non-identification means that if we have a $\pi(\mathbf{P}, \mathbf{G}, \mathbf{H}, \theta_1, \dots, \theta_S)$ in (3.5) that is generated from one of these points, then we will not be able to recover a unique value of $(\mathbf{P}, \mathbf{G}, \mathbf{H}, \theta_1, \dots, \theta_S)$

from those π 's. The most interesting cases in terms of its economic interpretation are straight forward to generalize to higher T and S . Some of the other points we have located are only for that simpler case. These other cases usually do not have an economic interpretation in our model but impose very particular restrictions between the different points of support of the unobserved heterogeneity.⁵ In practical terms we are not generally interested in such cases. Aside for the obvious cases where any of the parameters is at the boundaries of the parameters' space (that is, it is zero or unity), the following is a list of the interesting non-regular points we have found⁶:

1. $P_s = G_s = H_s$ for all $s = 1, \dots, S$. In this case the model is not a Markov Chain but a static model where each period, including the initial observation, are independent realizations of mixtures of identical Bernoulli distributions. If we generate the π from this case for specific values and then invert the system (3.5) to recover $(\mathbf{P}, \mathbf{G}, \mathbf{H}, \theta_1, \dots, \theta_S)$, we recover a line of values (of measure zero) where one of the parameters is left unidentified and the other parameters are a function of it. So in this case, (3.9) is necessary but not sufficient for identification.
2. P_s is from the steady state. That is: $P_s = \frac{G_s}{1+G_s-H_s}$, for all $s = 1, \dots, S$. As said for all the cases in this list, we have found that when $T = \min T$ not only the Jacobian does not have full rank, but also it is not locally identified. Of course if we know that our initial observation is from the steady state and incorporate this to the model we try to identify, then (3.9) is again a sufficient condition for identification.
3. $G_1 = G_2 = \dots = G_S$, or $H_1 = H_2 = \dots = H_S$, or $P_1 = P_2 = \dots = P_S$. In this case, the S points of support are not distinct points in all the three dimensions. This is a violation of the assumption that the model we tried to identified has S points of support. In practice, this is the non-identified case with probably easier solution because we only need to adjust S to the actual (smaller) number of distinct points of support that define the model.

3.2.2. Global Identification

Even if we are usually interested on global identification, the previous results are still useful because local identification is necessary for global identification. However, it is not sufficient in general and we still may want to obtain conditions that guarantee

⁵An example of this are non-regular points that require $G_1 = (1 - H_2)$, $G_2 = (1 - H_1)$, plus other restrictions on the P 's and θ .

⁶Regardless of the effort made to cover all possibilities (at least in the $T = 2$, $S = 2$ case), we can not show that this list is exhaustive.

global identification. The problem is that, as explained in Rothenberg (1971), it is much more difficult to prove, and there are few global identification results.

We tried to show global identification starting with the simpler case with $T = 2$, $S = 2$, but we have not been successful. In that case, a computer program using symbolic calculus is able to invert (3.5) when giving specific number to π . Based on doing this for many values of the π we conjecture that all the locally identified points are also globally identified. However proving this conjecture is correct remains an open question. This means the only global identification results we have for this model are those in Browning and Carro (2011) that, as explained, are not based on all the possible different moment conditions.

4. Estimation and Testing against alternative models.

4.1. The time-homogeneous first order Markov model.

The identification analysis above suggests the following estimation procedure. First, estimate the proportions for each path and test for the model restrictions. If these are not rejected, then impose the conditions and solve for the unknown parameters using the identification conditions.

In practice, it is much better and more efficient to combine the two steps in a maximum likelihood analysis. This is particularly the case given we cannot derive analytically the inequality constraints that the HFOM imposes (see the discussion in subsection 2.4).

Take the full heterogeneity model with $S = \Upsilon_T$ so that we have a just identified model. Using the first form in (3.3), the structural model is:

$$\pi_j = \sum_{s=1}^S p_j (P_s, G_s, H_s) \theta_s \quad j = 1, 2, \dots, \Gamma \quad (4.1)$$

Define a indicator $\delta_{ij} = 1$ if unit i has path j and zero otherwise. For given parameters, the likelihood of a sample $\{y_{i0}, y_{i1}, \dots, y_{iT}\}_{i=1,2,\dots,N}$ is:

$$\prod_{i=1}^N \prod_{j=1}^{\Gamma} \left(\sum_{s=1}^S p_j (P_s, G_s, H_s) \theta_s \right)^{\delta_{ij}} = \prod_{j=1}^{\Gamma} \left(\sum_{s=1}^S p_j (P_s, G_s, H_s) \theta_s \right)^{n_j} \quad (4.2)$$

where n_j is the number of times a sequence j appears in the sample. Denote the sample proportions for path j $c_j = n_j/N$. The log-likelihood function for the mixture

model is:

$$\ell_{mix} = \sum_{i=1}^N \sum_{j=1}^{\Gamma} \delta_{ij} \log \left(\sum_{s=1}^S p_j (P_s, G_s, H_s) \theta_s \right) \quad (4.3)$$

$$= N \sum_{j=1}^{\Gamma} c_j \log \left(\sum_{s=1}^S p_j (P_s, G_s, H_s) \theta_s \right) \quad (4.4)$$

Note that N is irrelevant for the maximization. With an iid random sample $c_j \rightarrow \pi_j$ as $N \rightarrow \infty$. from the structural model. The advantage of using the likelihood framework for estimation is that we know how to use all the information on the sample, how to make inference and how to test different models.

4.2. The unrestricted model.

A natural benchmark against which to test the HFOM model is the saturated model with:

$$\begin{aligned} S &= \Gamma, \mathbf{A} = I, \theta = \pi \text{ or} \\ S &= 1, \mathbf{A} = \pi \end{aligned} \quad (4.5)$$

These both give the likelihood value:

$$\ell_{sat} = \sum_{i=1}^N \sum_{j=1}^{\Gamma} \delta_{ij} \log (c_j) = N \sum_{j=1}^{\Gamma} c_j \log (c_j). \quad (4.6)$$

This can be used to derive a likelihood ratio statistic for the test of the Markov model against the unrestricted alternative. In particular, if we do not reject the restriction from (4.6) to (4.4) then we cannot reject that we have a time-homogeneous first order model. In practice, the large number of zeros for most paths if T is moderately sized leads to a distribution for the LR statistic that is very far from a χ^2 distribution with degrees of freedom equal to the number of restrictions (R_T in table 2.2). In this case, we should simulate the distribution of the LR statistic to calculate the correct probability of the observed LR statistic.

4.3. The unrestricted HFOM model or restricted saturated model.

We can also write a closed form expression for the model with the HFOM equality restrictions from subsection 2.4 imposed, using equation (2.15). Let $k(j)$ denote the group (running from $k = 1, ..r_T$) that path j belongs to. Then define predicted

probabilities for path $j = 1, \dots, \Gamma$ by:

$$\hat{c}_j = \frac{1}{n_{k(j)}} \sum_{j \in k(j)} c_j \quad (4.7)$$

That is, we replace the unrestricted proportions for each path by the mean for the group.⁷ The likelihood function is then given by:

$$\ell_{res_sat} = \sum_{i=1}^N \sum_{j=1}^{\Gamma} \delta_{ij} \log(\hat{c}_j) = N \sum_{j=1}^{\Gamma} c_j \log(\hat{c}_j) \quad (4.8)$$

This likelihood function also plays an important role in the estimation and choice of the mixing model. If we take a mixture with the maximal number of components, Υ_T in Table 3.1 then it has a log likelihood value that is bounded above by ℓ_{res_sat} . The mixture model will only attain this likelihood value if the observed $\hat{\mathbf{c}}$ vector satisfies the inequality constraints discussed in subsection 2.4. Given the difficulties of finding global maxima when we have many components, having a benchmark value is a considerable advantage. Denote the likelihood value of this mixture model by ℓ_{mix}^{Υ} . Now consider a model with fewer than the maximum number of points of support: $S < \Upsilon_T$. We have the following ordering for the likelihood function values:

$$\ell_{sat} \geq \ell_{res_sat} \geq \ell_{mix}^{\Upsilon} \geq \ell_{mix}^S \quad (4.9)$$

As already discussed, the likelihood ratio statistic does not have a known general distribution (see chapter 6.4 of McLachlan and Peel (2004)) but a test of the model with a smaller number of points of support than Υ_T can be constructed based on the simulated distribution for the LR statistic, taking the restricted model as the null.⁸

If we reject the first order time-homogeneous model, we have a number of alternatives. We could try a time-homogeneous second order model; this would give rise to similar calculations to those made above. Alternatively, we could continue to maintain that the model is a first order Markov chain but with time-inhomogeneous transition probabilities. One variant would be to assume a structural break. A second variant has that the transition probabilities depend on observable time-varying covariates. We consider that in the next section.

⁷To illustrate, consider the case $T = 3$. Paths 3 (0010) and 5 (0100) are restricted in the HFOM model to have the same probability and so are paths 12 and 14. Therefore, $\hat{c}_3 = \hat{c}_5 = \frac{c_3 + c_5}{2}$; $\hat{c}_{12} = \hat{c}_{14} = \frac{c_{12} + c_{14}}{2}$; $\hat{c}_j = c_j$, all other j .

⁸Given that we have a fully parametric model, simulating the distribution of the LR statistic under the null seems preferable to subsampling methods.

4.4. Testing for a second order Markov process

Although the test of the HFOM model against the saturated model allows for any alternative, it may lack power since the alternative is not specified. The obvious alternative is a time-homogeneous second order process. Given the estimates of the first order process, we can derive a standard LM test for this. The log-likelihood of a time-homogeneous second order Markov process has the following form for the predicted probabilities:

$$\begin{aligned}
p_j (P_{00s}, P_{01s}, P_{10s}, G_{00s}, G_{10s}, H_{01s}, H_{11s}) = & \\
& P_{00s}^{1(y_0^j=0, y_1^j=0)} P_{01s}^{1(y_0^j=0, y_1^j=1)} P_{10s}^{1(y_0^j=1, y_1^j=0)} * \\
(1 - P_{00s} - P_{01s} - P_{10s})^{1(y_0^j=1, y_1^j=1)} G_{00}^{n_{00}^j} (1 - G_{00})^{n_{00}^j} G_{10}^{n_{10}^j} * & \\
(1 - G_{10})^{n_{10}^j} H_{01}^{n_{01}^j} (1 - H_{01})^{n_{01}^j} H_{11}^{n_{11}^j} (1 - H_{11})^{n_{11}^j} & \quad (4.10)
\end{aligned}$$

where $1(\cdot)$ is the indicator function and:

$$P_{01} = \Pr(y_{i0} = 0, y_{i1} = 1), \quad (4.11)$$

$$G_{10} = \Pr(y_{it} = 1 \mid y_{it-2} = 1, y_{it-1} = 0), \quad (4.12)$$

$$H_{01} = \Pr(y_{it} = 1 \mid y_{it-2} = 0, y_{it-1} = 1), \quad (4.13)$$

...

This has seven parameters per type s , instead of three. Three of them are to account for the initial conditions, since now we have to condition on two previous observations. The other four are the transition probabilities given by the second order Markov process, what imposes less restrictions on the data than the first order process. Therefore, the log-likelihood now depends on $8S - 1$ parameters.⁹ To perform the LM test we have to:

1. Derive the log-likelihood with respect to $\{P_{00s}, P_{01s}, P_{10s}, G_{00s}, G_{10s}, H_{01s}, H_{11s}\}_{s=1}^S$ and $\{\theta_s\}_{s=1}^{S-1}$. This gives the score vector denoted by $g(\cdot)$, and allows us to calculate the outer-product of the score, denoted by $h(\cdot)$.
2. Evaluate $g(\cdot)$ and $h(\cdot)$ at the estimated values of the parameters of the first order Markov model $\left(\left\{ \hat{P}_s, \hat{G}_s, \hat{H}_s \right\}_{s=1}^S, \left\{ \hat{\theta}_s \right\}_{s=1}^{S-1} \right)$. This means that we evaluate $g(\cdot)$ and $h(\cdot)$ at $P_{00s} = (1 - \hat{P}_s) (1 - \hat{G}_s)$, $P_{01s} = (1 - \hat{P}_s) \hat{G}_s$, $P_{10s} =$

⁹This means that we are keeping S constant. Related with this, it is important to notice that to point identify a first order Markov model with S points of support does not imply that a second order Markov model with S points of support can also be point identified.

$\widehat{P}_s(1 - \widehat{H}_s)$, $G_{00s} = G_{10s} = \widehat{G}_s$, $H_{01s} = H_{11s} = \widehat{H}_s$ for $s = 1, \dots, S$, and $\{\widehat{\theta}_s\}_{s=1}^{S-1}$. Denote the values we get from this by \widehat{g} and \widehat{h} .

3. Then, the test statistic is

$$LM = \widehat{g}'\widehat{h}^{-1}\widehat{g} \quad (4.14)$$

Under the standard regularity conditions this test statistic is asymptotically distributed as χ_b^2 . The degrees of freedom are

$$b = (7S + S - 1) - (3S + S - 1) = 4S \quad (4.15)$$

4.5. Homogenous marginal dynamic effect.

We shall not consider the homogeneous case with (G, H, P) the same for everyone, since it is hardly considered a possibility. A less restricted model than the homogeneous case is the usual ‘fixed effect’ case which only allows for one source of unobservable heterogeneity. The latter is usually in the intercept of the index in (1.1). A close analogue here is that we have a homogeneous dynamic marginal effect:

$$H_i = M + G_i \text{ for some constant } M \in [-1, 1] \quad (4.16)$$

This test can be done using a standard LR test statistic of the $(S - 1)$ restrictions imposed.

4.6. Testing for time homogeneity.

As well as testing against a specific time homogeneous model, we can also derive a test for time homogeneity. To do this, we split the sample into an estimation sample $\{y_{i0}, y_{i1}, \dots, y_{iE}\}$ and a hold-out sample $\{y_{iE+1}, y_{iE+2}, \dots, y_{iT}\}$. We estimate the mixture model on the estimation subsample and test whether the predictions for the hold-out subsample fit. To do this we take the same transition probabilities for the hold-out subsample. To generate the distribution for period $E + 1$ (the initial period for the hold-out sample) we use the estimated probabilities and the Chapman-Kolmogorov equations to generate the relevant distribution. An alternative procedure is split the sample into two equal subsamples in terms of term $(E \text{ close to } (T + 1) / 2)$, estimate on each subsample separately and then test whether the two sets of estimates are statistically different. A particularly simple variant of a stability test of this sort this will be given in the empirical section.

5. Allowing for covariates.

5.1. Model and Parameters of Interest

In the presence of covariates in the model, our estimation is conditional on the covariates, which are assumed to be strictly exogenous. As before, we look directly at the conditional probabilities:

$$\begin{aligned} H_{xi} &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 1, x_{it} = x) \\ G_{xi} &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 0, x_{it} = x) \end{aligned}$$

where H_{xi} and G_{xi} are defined for each value x of x_{it} , and at the unconditional probability of a unit value for the initial observation:

$$P_{xi} = \Pr(y_{i0} = 1 \mid x_{i0} = x) \quad (5.1)$$

In addition to the marginal dynamic effect and the long run proportion of unit values mentioned for the model without covariates, (P_{xi}, G_{xi}, H_{xi}) can be used to generate any other outcomes or parameters of interest. There are several candidates but the most widely considered are those informing about the *marginal effects* of the change in a explanatory variable x over the probability of y_{it} being equal to 1, (considering only one covariate for notational convenience):

$$\begin{aligned} M_{x'i} &= \Pr(y_{it} = 1 \mid y_{i,t-1}, x_{it} = x'' = x' + 1) - \Pr(y_{it} = 1 \mid y_{i,t-1}, x_{it} = x') \\ &= \begin{cases} H_{x''i} - H_{x'i} & \text{if } y_{it-1} = 1 \\ G_{x''i} - G_{x'i} & \text{if } y_{it-1} = 0 \end{cases} \end{aligned} \quad (5.2)$$

Given that the marginal effects are heterogenous across individuals in the population, the interest is usually in knowing their distribution or some moments. There are many possible measures that could be considered. For example, there is more than one mean effect that could be of interest. Here we mention just two of them.

The first is expected effect on the probability of $y = 1$ of a change in variable x given the distribution of the unobservables conditional on $x = x'$:

$$\begin{aligned} E_{(G,H)|x} [M_i | x'] &= \int \int [(H_{x''i} - H_{x'i}) \Pr(y_{it-1} = 1 | x_{it} = x'; G_{xi}, H_{xi}) \\ &+ (G_{x''i} - G_{x'i}) \Pr(y_{it-1} = 0 | x_{it} = x'; G_{xi}, H_{xi})] dF_{(G,H)|x}(G_{xi}, H_{xi} | x') \end{aligned} \quad (5.3)$$

This is equivalent to the parameter of interest estimated in Altonji and Matzkin (2005). If x were a treatment indicator variable with $x' = 0$ and $x'' = 1$, then (5.3)

would give the Average Treatment on the Untreated effect.

The second example is the average marginal effect without conditioning on x :

$$\begin{aligned}
E_{(G,H)} [M_i] &= \int \int [(H_{x''i} - H_{x'i}) \Pr(y_{it-1} = 1 | G_{xi}, H_{xi}) \\
&+ (G_{x''i} - G_{x'i}) \Pr(y_{it-1} = 0 | G_{xi}, H_{xi})] dF_{(G,H)}(G_{xi}, H_{xi})
\end{aligned} \tag{5.4}$$

This is equivalent to the parameter of interest proposed by Chamberlain (1984) defined there as the expected effect for a randomly drawn individual. If x were a treatment indicator variable with $x' = 0$ and $x'' = 1$, then (5.4) would give the Average Treatment Effect.

Equations (5.3) and (5.4) are answers to different questions and with more explanatory variables, averages over different distributions could be considered. On the other hand, as explained in the introduction, average effects may not be very informative in nonlinear models such as this. In such a case other moments of the individual marginal effects like the median are more informative. Furthermore, there are cases where the entire distribution of the marginal effect over the population is the object of interest; see Heckman, Smith and Clements (1997). In the IO literature the object of interest is the entire distribution of the individual price elasticities; see, for example, Nevo (2001). Identifying and estimating the distribution of (P_i, G_{xi}, H_{xi}) allow us to obtain any possible parameter of interest since it fully characterizes the HFOM model.¹⁰

5.2. Identification with exogenous covariates

Adding covariates not only changes the number of transition probabilities we have to identify, but also introduces the possibility of dependence between the probability of being of each unobserved type and the covariates. We begin with the simplest case, a binary covariate that is constant over time, as an introduction to the case with a general x_{it} . The special case of covariates that only vary with time (that is, in each t they take a common value for all i) is explicitly discussed, including the case with time dummies. A summarizing table with numbers for representative cases can be found at the end of this section.

If we only have an x variable that is constant over time and only varies across individuals (for example, year of birth or education), it is straightforward to extend our identification result in section 3.2. For a binary x_i , the time homogeneous first

¹⁰Since the x variables are assumed to be exogenous, there is no problem in obtaining their distribution from a random sample when needed.

order Markov model is fully characterized by:

$$\begin{aligned}
P_{0i} &= \Pr(y_{i0} = 1 \mid x_i = 0); P_{1i} = \Pr(y_{i0} = 1 \mid x_i = 1) \\
G_{0i} &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 0, x_i = 0); H_{0i} = \Pr(y_{it} = 1 \mid y_{i,t-1} = 1, x_i = 0) \\
G_{1i} &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 0, x_i = 1); H_{1i} = \Pr(y_{it} = 1 \mid y_{i,t-1} = 1, x_i = 1) \quad (5.5)
\end{aligned}$$

As before, we consider a nonparametric discrete distribution for $(P_{0i}, G_{0i}, H_{0i}, P_{1i}, G_{1i}, H_{1i})$ with S distinct points of support $\{P_{0s}, G_{0s}, H_{0s}, P_{1s}, G_{1s}, H_{1s}\}_{s=1}^S$ with probabilities given by the $(S \times 1)$ vector θ_x with non-negative values θ_{xs} that sum to unity.

$$\theta_x = \begin{cases} (\theta_{01}, \dots, \theta_{0S})' & \text{if } x = 0 \\ (\theta_{11}, \dots, \theta_{1S})' & \text{if } x = 1 \end{cases} \quad (5.6)$$

where each vector sum one.¹¹

Here we simply divide the observation in two groups (one with $x_i = 0$ and the other $x_i = 1$) and do the identification analysis and estimation for each group. Each group contains the same number of parameters to identify and the same moment conditions as the problem without covariates. Therefore the number of period necessary and sufficient to identify $(P_{0i}, G_{0i}, H_{0i}, P_{1i}, G_{1i}, H_{1i})$ with S points of support almost everywhere is the same as to identify (P_i, G_i, H_i) with S points of support in the case without covariates in section 3.2. If x_i takes N_x values then we stratify based on the value of x_i and everything is the same as with a binary x_i .

So far, adding an x_i has not changed the periods needed to identify S points of support. A case in which identification can be improved is if the probability of each type is independent of x_i . That is $\theta_x = (\theta_1, \dots, \theta_S)'$ for all values of x . This reduces the number of parameters, but not the number of equations. There are $(3SN_x + (S - 1))$ parameters instead of $(4N_x(S - 1))$. Therefore, to identify S points of support in this case we need $S \leq \frac{N_x r_T - N_x + 1}{3N_x + 1}$.

We now consider the case of x_{it} covariates that have positive probability of taking any value of their support at any i and t . Maintaining the independent assumption

¹¹The analysis and estimation is made conditional on X , and therefore we are specifying and obtaining the distribution of the individual parameters conditional on x . Nevertheless, the unconditional distribution can be calculated from this conditional distribution and the distribution of x , which can be obtained from the data.

T	2	3	4	5	6	7	8
$r_{xit}(T, 2)$	60	184	472	1056	2132	3976	6964
$r_{xit}(T, 4)$	464	2656	12088	45888	151456	447648	1210032
$r_{xit}(T, 6)$	1548	12984	84852	454104	2079840		

Table 5.1: Number of independent paths. Discrete covariate

$\Pr(s|x_{i0}, \dots, x_{iT}) = \Pr(s)$, for each point of support s :

$$\begin{aligned}
P_{xs} &= \Pr(y_{i0} = 1 \mid x_{i0}, s) \\
G_{xs} &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 0, x_{it}, s) \\
H_{xs} &= \Pr(y_{it} = 1 \mid y_{i,t-1} = 1, x_{it}, s) \\
\theta_s &= \Pr(s|x_{i0}, \dots, x_{iT}) = \Pr(s)
\end{aligned} \tag{5.7}$$

It is conceptually simple to extend our model if the additional covariates are discrete. We denote the number of values that a discrete x_{it} can take by N_x . The number of parameters to identify is $(3N_x + 1)S - 1$. The probability of a path j given $\{x_{it}\}_{t=1}^T$ and (P_{xs}, G_{xs}, H_{xs}) is:

$$p_{js} = P_{x_{i0}s}^{y_0^j} (1 - P_{x_{i0}s})^{(1-y_0^j)} \prod_x G_{xs}^{n_{01|x}^j} (1 - G_{xs})^{n_{00|x}^j} H_{xs}^{n_{11|x}^j} (1 - H_{xs})^{n_{10|x}^j} \tag{5.8}$$

where \prod_x denotes the product over the N_x values of x_{it} , $n_{01|x}^j$ is the number of $y_{it-1} = 0 \rightarrow y_{it} = 1$ transitions given $x_{it} = x$ for path j , . The number of possible paths in our system is $2^{2(T+1)}$, because we have 2^{T+1} possible paths of $\{y_{it}\}_{t=0}^{T+1}$ given each one of the 2^{T+1} possible observations of $\{x_{it}\}_{t=0}^{T+1}$. As in other cases, some of those paths will give the same equation. We denote the number of different equations by $r_{xit}(T, N_x)$ whose specific expression is in equation (A.2), given and proved in Appendix A. Table 5.1 shows this number for some T and N_x . Notice that $r_{xit}(T, N_x)$ grows very fast with N_x .

By the same arguments used in the case without covariates, the number of periods needed to identify S points of support almost everywhere is given by the following condition

$$S \leq \frac{r_{xit}(T, N_x) - N_x^{T+1} + 1}{3N_x + 1} \tag{5.9}$$

obtained from comparing the number of parameters with the rank of the Jacobian of the system. Table 5.2 gives in each column the highest value of S for which T in that column is $\min T$ for identification of that S (that is, the maximum integer value of S such that (5.9) is satisfied for each T).

Continuous covariate If x_{it} is a continuous variable instead of discrete, we discretise it on an arbitrary number N_x of grid points and use the previous result. Since (5.9) is increasing with N_x , this means we can potentially nonparametrically identify as many points of support as we wish simply by discretising the continuous covariate in as many points as needed; see Remark 2 (iv) in Kasahara and Shimotsu (2009).

5.2.1. Semiparametric model

In the previous analysis we have not only allowed for maximal (nonparametric discrete) heterogeneity across i , but also we are not restricting our HFOM model to have a particular functional form. In particular we have not imposed any restriction on the way different values of x_{it} affects y_{it} . Nevertheless, if x_{it} is continuous, or a cardinal discrete variable that takes many values, such as year of birth, then the effect of different values of x is usually restricted by a parametric form. The obvious example is a linear index model:

$$\begin{aligned}
 P_{si} &= F_0(p_{s0} + p_{s1}x_{i0}) \\
 G_{sit} &= F(g_{s0} + g_{s1}x_{it}) \\
 H_{sit} &= F(h_{s0} + h_{s1}x_{it}) \\
 \theta_{si} &= F_\theta(d_{s0} + \sum_{t=0}^T d_{s1t}x_{it})
 \end{aligned} \tag{5.10}$$

F_0 , F and F_θ are known cdf functions, like the standard normal cdf or the standard logistic function. This is equivalent to the representation

$$\Pr(y_{it} = 1 \mid y_{i,t-1}, x_{it}) = F(\eta_i + \alpha_i y_{i,t-1} + \beta_i x_{it} + \delta_i x_{it} y_{i,t-1})$$

where $(\eta_i, \alpha_i, \beta_i, \delta_i)$ follow a discrete distribution with S points of support. Equation (5.10) allows for dependence between θ and x , and includes the independent case ($d_{s1t} = 0$ for $t = 0, \dots, T$) and a case with correlation only with the observation of the initial period ($d_{s1t} = 0$ for $t = 1, \dots, T$).

The number of parameters is now $(8 + T)S - (T + 2)$. It does not depend on the number of values x_{it} can take. This is reducing the number of parameters to identify with respect to to the non-parametric case without altering the number of different moment conditions, $r_{xit}(T, N_x)$. The latter number still depends on N_x and it is given by equation (A.2). This implies that the maximum number of points of support for

which T periods are required for identification is

$$\frac{r_{xit}(T, N_x) - N_x^{T+1} + T + 2}{8 + T} \quad (5.11)$$

Assuming independence between θ and x , this number is $\frac{r_{xit}(T, N_x)+1}{7}$, which is clearly smaller than (5.9). This reflects the important gains of the semiparametric assumption.

5.2.2. Time dummies and Common variables

Finally, we consider the situation in which we add a covariate that it is common to all individuals and only varies across periods: $x_{it} = x_t$ for all i . For instance, this is the case with aggregate variables being used in a micro study, or with time dummy variables. Since we are studying identification over the i population for a fixed T , we are only going to observe a given and fixed realization of $\{x_t\}_{t=1}^T$. This implies we only have the 2^{T+1} possible paths given $\{x_t\}_{t=1}^T$ that arises from the possible combinations of $\{y_{it}\}_{t=1}^T$ we can observe over the population of i . Then, the number of equations in our system here are the same as in the case without covariates and the rank of the Jacobian also depends on r_T . For the same reason, x_t is not going to be an informative variable for the probability of y_{i0} , nor for the distribution of the heterogenous parameters over the i population, that is, $\Pr(s | \{x_t\}_{t=1}^T) = \Pr(s) = \theta_s$. However, this covariate increases the number of parameters to be identified. Therefore, this is the only case where more periods are required for identification than in the case without covariates

A situation often found in practice is the use of time dummies. These variables take deterministic values, and, while treated as separate variables, the only meaningful situation is where one of them takes value one and all the other take value zero. If we add time dummies to the model, we have $K = T$ variables x_t that can take $N_x = 2$ values each, but in a deterministic way. Thus we have $(2 + 2T)S - 1$ parameters: one G and H for each time dummy. Then,

$$S \leq \frac{r_T}{2 + 2T} \quad (5.12)$$

This implies a much larger T to identify a given S . For example, we need $T \geq 8$ for the identification conditions of a model with $S = 4$ to be satisfied. $T = 23$ is the minimum T we need to identify $S = 11$.

If, on the other hand, \mathbf{X}_t contains K discrete variables taking many values or continuous variables, then we can use a semiparametric model to capture the effect

T	2	3	4	5	6	7	23
Υ_T : No covariates	2	3	5	8	11	14	138
Covariate constant over time ($x_{it} = x_i$ for all t)							
Any N_x , free relation with θ	2	3	5	8	11	14	138
$N_x = 10$, independence of θ	2	4	6	10	13	18	178
$N_x = 10$, semiparametric	9	16	26	39	54	71	691
Covariates $x_{it} = x_t$ for all i							
Time dummies	1	1	2	2	3	3	11
2 continuous x_t , semiparametric	1	1	2	4	5	7	69
Covariate that varies in both i and t							
$N_x = 2$, independent of θ	7	24	63	141	286	531	
$N_x = 4$, independent of θ	30	184	851	3214	10390	29393	
$N_x = 4$, semiparametric	40	218	992	3215	9648	25474	
$N_x = 6$, semiparametric	133	1063	6423	31342	128565		
N_x is the number of possible values x can take. Where semiparametrically is not specifically mentioned, a nonparametric first HFOM model with the indicated covariates is being considered.							

Table 5.2: Maximum number of points of support for some representative cases

of X . For each point of support s :

$$\begin{aligned}
 G_s &= F(g_{s0} + \sum_{k=1}^K g_{sk} x_{kt}) \\
 H_s &= F(h_{s0} + \sum_{k=1}^K h_{sk} x_{kt})
 \end{aligned}
 \tag{5.13}$$

where F is a known cdf, such as the logistic. In this case the number of parameters is $(2 + 2(K + 1))S - 1$, and

$$S \leq \frac{r_T}{2 + 2(K + 1)}
 \tag{5.14}$$

For example, if $K = 2$ and $S = 9$, then $\min T = 8$; or if $K = 2$ and $S = 69$, then $\min T = 23$. This and values for other cases can be found in table 5.2.

6. An empirical illustration.

6.1. Sample selection.

We consider the incidence of unemployment in a year for workers in Denmark from 1980 to 2003 (so that $T = 23$). We draw a sample of male workers with high school education who were aged 25 at the beginning of 1980 and who are continuously married to the same wife for all 24 years that we follow them. This is thus a *very*

homogeneous sample in terms of observables; we do this so that our finding of considerable heterogeneity cannot be attributed to insufficient allowance for observable heterogeneity. In all, we have 2571 such workers.¹² We create a dummy variable y_{it} which is set to unity if worker i has any unemployment in year t (and zero otherwise). The following Table gives some statistics for the sample.

	Number	Proportion
Total sample size	2571	—
No unemployment	936	36.4
At most 1 year with unemployment	1141	44.4
At most 2 years with unemployment	1291	50.2
At most 3 years with unemployment	1435	55.8
At most 5 years with unemployment	1710	66.5
At most 10 years with unemployment	2188	85.1
At most 20 years with unemployment	2519	98.0
Unemployment in all years	16	0.6

Table 6.1: Incidence of unemployment

6.2. The model without covariates.

The indicator variable y_{it} is unity if worker i had a spell of unemployment in year t . We begin with the model without covariates. The likelihood function value for the saturated model, ℓ_{sat} (4.6), is $-12,252$. The value for the saturated HFOM model, ℓ_{res_sat} , (4.8), is $-17,449$. The likelihood ratio statistic, $2(\ell_{sat} - \ell_{res_sat})$, is thus $10,395$.¹³ When estimating the mixture model we restrict the mixing probabilities $\theta_s \geq 0.01$ and we restrict G_s , H_s and P_s to be between 0.01 and 0.99 to ensure that we do not assign zero probability to any path. The maximum number of support points we could have for the HFOM model is 138 (see Table 3.1). In practice, we cannot find more than a much smaller number than this; see Table 6.2. For ease of reading, we present all likelihood function values for mixture models in LR terms relative to the value for ℓ_{res_sat} ; that is, the LR statistic shown is $2(\ell_{res_sat} - \ell_{mix}^S)$. We also show how many mixing parameters are at the imposed minimum of 0.01. As can be seen, it does not seem to be possible to estimate with more than nine components; that is, $\ell_{mix}^{10} \simeq \ell_{mix}^9$.

¹²Denmark has an administrative panel that follows *all* of the population of about five million from 1980 onwards. Consequently we can select very homogeneous strata without compromising sample size. Indeed, the sample drawn here is, in fact, the population of men who fulfilled the selection criteria.

¹³In an earlier version of this paper we developed a parametric bootstrap test for assessing whether the HFOM hypothesis is rejected and for choosing S if it is not. Since this is controversial (see Feng and McCulloch (1996)) and takes us too far from the main theme of this paper, we do not present results here. In the next section we develop a valid test against an *HFOM* with covariates.

S	df	LR stat	$\# \theta'_s = 0.01$
2	547	1,063	0
3	543	701	0
4	539	605	0
5	535	536	0
6	531	512	0
7	527	500	0
8	523	494	0
9	519	491	1
10	515	491	2

Table 6.2: Fit for different numbers of support points

Since we are concerned to illustrate the mechanics of our method, we shall sidestep the issue of the distribution of the LR statistics and simply take a convenient value, $S = 5$. Table 6.3 presents the estimates for the model with 5 points of support. These display a number of features. First, all groups display positive state dependence ($H_s > G_s$). Second, the marginal dynamic effects ($H_s - G_s$) vary quite considerably across groups. The LR statistic for the hypothesis of a homogeneous marginal dynamic effect,

$$H_s = G_s + (H_1 - G_1) \text{ for } s = 2, \dots, 5 \quad (6.1)$$

is 421; this is distributed as a $\chi^2(4)$ and represents a decisive rejection of this homogeneity assumption. Moreover the (weighted) correlation between G and H is -0.35 ; the conventional ‘one fixed effect’ assumption imposes that the correlation is positive so that even the qualitative implication is wrong for the homogeneous model.

Group	Probabilities				
	P	G	H	M	θ
	$p(y_0 = U)$	$p(U E)$	$p(U U)$	$H - G$	Proportion
1	0.27	0.01	0.87	0.86	0.34
2	0.64	0.10	0.69	0.59	0.28
3	0.01	0.03	0.48	0.46	0.24
4	0.73	0.36	0.82	0.46	0.08
5	0.25	0.18	0.34	0.16	0.06

Table 6.3: Parameter estimates with five support points

To see the substantive implications of the estimates it is best to graph the implied paths for the probability of being unemployed at some time during the year. This is shown in the left panel of Figure 6.1 which graphs the probabilities implied by the Chapman-Kolmogorov equations for the five groups against age (or year, since all the workers in the sample are in the same birth cohort). The groups can be identified

from their initial values given in Table 6.3. The figure suggests a fascinating mix of workers who rarely experience unemployment (group 3), those who are very prone to unemployment (group 4) and those who start off badly, but quickly ‘find their feet’ (groups 2 and 1). However, there is evidence that the HFOM model does not fit the data well. This is shown in the right panel of the figure which shows the average proportions of unemployed for each year and the predicted mean from the model. The estimation imposes that the two coincide at age 25 but they are conspicuously different thereafter. A formal test for parameter stability can be constructed by splitting the sample and estimating with dummy shifters for H_s and G_s . If we do this with a dummy variable that is unity for the last 11 periods we have an LR statistic of 384; given that we have an extra parameter for each H_s and G_s this has a $\chi^2(10)$ distribution. This formally confirms the time inhomogeneity that we see in the right panel of Figure 6.1. To capture this time-inhomogeneity we turn to estimation adding the covariates to the model.

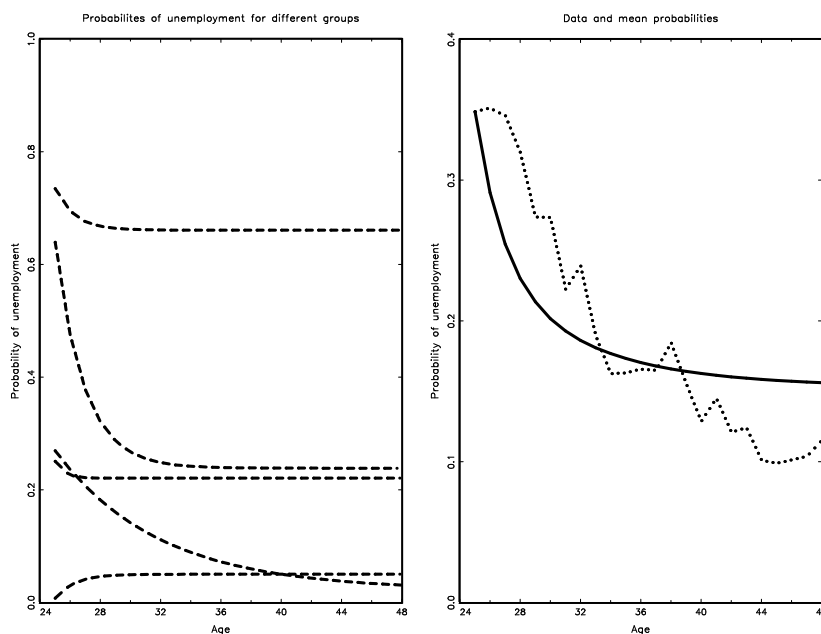


Figure 6.1: Probabilities with 5 points of support.

6.3. Model with covariates.

The right panel of Figure 6.1 suggests that we need to allow for time inhomogeneity that is associated with age. There also seem to be cyclical deviations from a smooth age profile. To capture these we include age and the aggregate unemployment rate

as covariates and the semiparametric specification in (??).^{14, 15} We continue to keep S fixed at 5. We first present likelihood ratio statistics for including the extra sets of variables. Since we have 5 points of support and we include regressors in the G_s and H_s transition probabilities, we have 10 extra parameters for each covariate. Table 6.4 presents the LR statistics against the model with 5 points of support and no covariates. As can be see, age and the aggregate unemployment rate are individually and jointly highly significant. Moreover, the $\chi^2(10)$ statistic for the stability test used in the previous subsection is 36; although formally this is a rejection, it is a considerable improvement on the model without age and cyclical effects.

Test against SFOM		
Model	df	χ^2
Age and cycle	20	808
Age only	10	766
Cycle only	10	163

Table 6.4: Tests for age and cyclical effects

As before, the implications of the estimates are most easily seen in figures of the unemployment sequences. These are given in figure 6.2. The right hand panel indicates that adding the age effects remedies most of the misfit seen in the earlier figure. The left hand panel shows that the impact of the business cycle is very heterogeneous. For example, the group who have very low probabilities are hardly affected at all. However, the next prone group (with a starting value of 0.22) display considerable cyclical variation. However, the group who have the highest propensity to be unemployed (the highest curve after age 32) also seem to be unaffected by the cycle. Thus the link between the propensity to be unemployed and the impact of the business cycle is not monotone. Estimates that did not allow for heterogeneity would mask this effect.

7. Conclusions.

This paper studies identification from a panel with given T of a non-parametric and a semiparametric dynamic binary choice model with maximal heterogeneity. The

¹⁴Note that aggregate unemployment rate is endogenous by definition, because the endogenous variable in our model is part of this explanatory variable. A solution to this is to construct an aggregate unemployment rate excluding from the population the group we are using. Since our group of workers represents less than 0.0001% of the working population, this will hardly have an impact on the estimates.

¹⁵Other factors that we could take into account are other macro variables such as changes in the UI system; individual time varying factors such as health or marital status and individual time invariant factors such as parental background. Note that in this empirical illustration we have taken account of the time invariant factor, cohort.

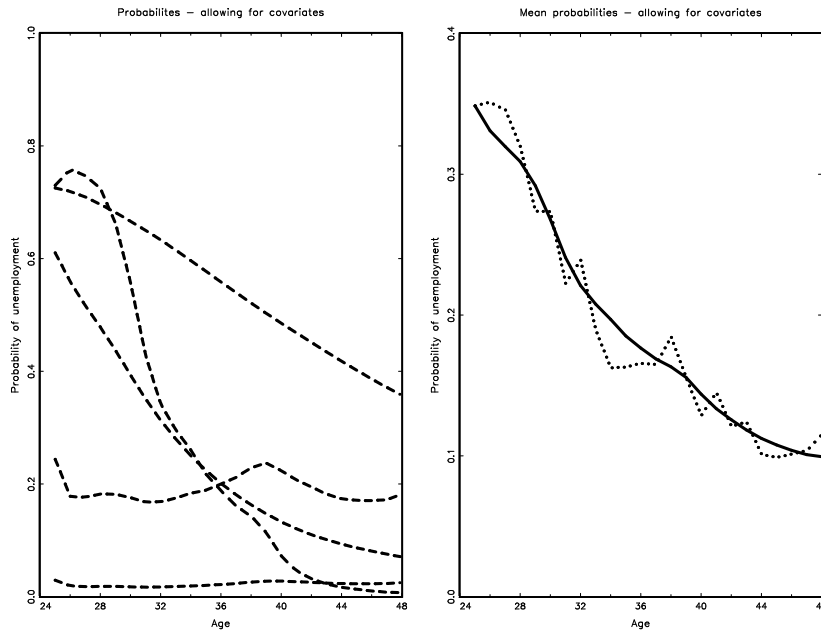


Figure 6.2: Probabilities with age and cyclical effects.

more traditional linear-index specification where only the constant term is individual specific is extended since the latter imposes undesired restrictions on the economic model and it does not generally fit the data. In contrast, our model allows variation in all of the parameters (and even the distribution function) across individuals. These models are not generally identified from a cross section of fixed- T periods.

In our specification the joint distribution of the initial observation and the transition probabilities is unrestricted, using nonparametric discrete mixture distributions. We establish necessary and sufficient conditions for point identification of our heterogeneity structure and show how it depends on the length of the panel.

A conclusion from this study is that a model with a very flexible distribution of the heterogeneity can be identified from a cross section of T periods, even for T as small as 3. The identification is strengthened if we have continuous covariates in the model. So a model that allows for maximal heterogeneity with a very rich and flexible distribution can be point identified. With such flexibility, important features of the distribution of the heterogeneity such as dependencies of transition probabilities on initial condition are unrestricted.

We show how to estimate using Maximum Likelihood. The asymptotic properties of the estimator in sample size with fixed panel length are well known: it is consistent and efficient. We apply the techniques we study to a long panel of Danish workers who are very homogeneous in terms of observables. One of our principal findings is that the impact of cyclical variations on unemployment for individual workers are

heterogeneous with non-obvious relations. Findings in this application seems to us very illustrative of the potential usefulness of our approach for applied work.

References

- [1] Alessie, R.; Hochguertel, S. and Soest, A. (2004): "Ownership of Stocks and Mutual Funds: A Panel Data Analysis." *Review of Economics and Statistics*, 86(3), pp. 783-96.
- [2] Altonji J. G. and R. L. Matzkin (2005): "Cross section and panel data estimators for nonseparable models with endogenous regressors", *Econometrica*, 73(4), 1053-1112.
- [3] Arellano, M. and B. H. Honoré (2001): "Panel Data Models: Some Recent Developments." *Handbook of Econometrics*, chapter 5, pp. 3229-96.
- [4] Becker, G. S.; Grossman, M. and Murphy, K. M. (1994) "An Empirical Analysis of Cigarette Addiction." *American Economic Review*, 84(3), pp. 396-418.
- [5] Bernard, A. B. and Jensen, J. B. (2004) "Why Some Firms Export." *Review of Economics and Statistics*, 86(2), pp. 561-69.
- [6] Browning, M. and J. M. Carro (2007). "Heterogeneity and Microeconometrics Modelling." *Advances in Economics and Econometrics, Theory and Applications: Ninth World Congress of the Econometric Society*, vol. 3. Cambridge University Press.
- [7] Browning, M. and J. M. Carro. (2010): "Heterogeneity in Dynamic Discrete Choice Models." *Econometrics Journal*, 13(1), pp. 1-39.
- [8] Browning, M. and J. M. Carro. (2011): "The identification of a mixture of first order binary Markov Chains" *unpublished manuscript*.
- [9] Carro J. M. and P. Mira (2006). "A dynamic model of contraceptive choice of Spanish couples", *Journal of Applied Econometrics*, 21, 955-980.
- [10] Chamberlain, G. (1984): "Panel Data", in Griliches, Z. and M.D. Intriligator (eds.) *Handbook of Econometrics*, vol. 2, Elsevier Science, Amsterdam.
- [11] Chernozhukov V., I. Fernandez-Val, J. Hahn and W. K. Newey (2009), "Identification and Estimation of Marginal Effects in Nonlinear Panel Data", *unpublished manuscript*.

- [12] Crawford G. S. and M. Shum (2005): "Uncertainty and Learning in Pharmaceutical Demand.", *Econometrica*, 73(4), 1137-1173.
- [13] Feng, Z. D. and C. E. McCulloch (1996). "Using Bootstrap Likelihood Methods in Finite Mixture Models", *Journal of the Royal Statistical Society, Series B*, 58(3), 609-617.
- [14] Fisher, F. M. (1966). *The Identification Problem in Econometrics*. New York. McGraw-Hill.
- [15] Gottschalk, P. and R. A. Moffitt (1994): "Welfare Dependence - Concepts, Measures, and Trends." *American Economic Review*, 84(2), pp. 38-42.
- [16] Ham, J. C. and L. Shore-Sheppard (2005): "The Effect of Medicaid Expansions for Low-Income Children on Medicaid Participation and Private Insurance Coverage: Evidence from the Sipp." *Journal of Public Economics*, 89(1), pp. 57-83.
- [17] Heckman, J. J. (1981) "Heterogeneity and State Dependence." *Studies in Labor Markets*, ch. 31, pp. 91-140.
- [18] Heckman, J. J. and Singer, B. (1984) "A Method for Minimizing the Impact of Distributional Assumptions in Econometric-Models for Duration Data." *Econometrica*, 52(2), pp. 271-320.
- [19] Heckman J. , J. Smith and N. Clements (1997):"Making the most out of Programme Evaluations and Social Experiments: Accounting For Heterogeneity in Programme Impacts", *Review of Economic Studies*, 64, 487-535.
- [20] Honorè, B. E. and E. Tamer (2006), "Bounds on Parameters in Panel Dynamic Discrete Choice Models", *Econometrica*, 74(3), pp. 611-629.
- [21] Hyslop, D. R. (1999): "State Dependence, Serial Correlation and Heterogeneity in Intertemporal Labor Force Participation of Married Women", *Econometrica*, 67, 1255-1294.
- [22] Kasahara, H. and K. Shimotsu (2009): "Nonparametric Identification of Finite Mixture Models of Dynamic Discrete Choices" *Econometrica*, 77(1), 135-175
- [23] Keane M. P. and K. I. Wolpin (1997): "The Career Decisions of Young Men", *Journal of Political Economy*, 105, 473-521.
- [24] McLachlan, G. and D. Peel (2004), *Finite Mixture Models*, Wiley-Interscience.

- [25] Nevo, A. (2001): “Measuring Market Power in the Ready-to-Eat Cereal Industry.”, *Econometrica*, 69(2), 307-342.
- [26] Rothenberg, T. J. (1971): “Identification in Parametric Models.”, *Econometrica*, 39(3), 577-591.

A. Proof of the number of different equations.

A.1. Number of ‘independent’ equations

Here we proof equation (2.14), that is, that the number of ‘independent’ equations in system (2.8) is

$$r_T = T(T + 1) + 2$$

By Lemma 2.1, all we have to do is to count the number of different sets $\{y_0^j, n_{00}^j, n_{01}^j, n_{10}^j, n_{11}^j\}$ that the $j = 1, \dots, 2^{T+1}$ possible paths can generate. Before counting, note that half of the r_T possible different paths have $y_0 = 0$ and the other half have $y_0 = 1$ and this two halves are symmetric, so we can count only paths with $y_0 = 0$ and multiply its number by two. Notice also that, for $y_0 = 0$ cases, $n_{00} + n_{01} > 0$, $n_{10} + n_{11} > 0$ only if $n_{01} > 0$, and that $n_{10} \in \{n_{01} - 1, n_{01}\}$. We set n_{00} to count, starting with the maximum value it can take:

- If $n_{00} = T$, then there is only one possibility: $\{(y_0, n_{00}, n_{01}, n_{10}, n_{11})\} = \{(0, T, 0, 0, 0)\}$
- If $n_{00} = T - 1$, then there is only 1 possibility: $\{(0, T - 1, 1, 0, 0)\}$
- If $n_{00} = T - 2$, then there are 2 possibilities: $\{(0, T - 2, 1, 1, 0), (0, T - 2, 1, 0, 1)\}$
- If $n_{00} = T - 3$, then there are 3 possibilities: $\{(0, T - 3, 2, 1, 0), (0, T - 3, 1, 1, 1), (0, T - 3, 1, 0, 2)\}$
- If $n_{00} = T - m$, then there are m possibilities, which are:

$$\left\{ \left(0, T - m, \left\lceil \frac{m - q}{2} \right\rceil, \left\lfloor \frac{m - q}{2} \right\rfloor, q \right) \right\}_{q=0}^{m-1} \quad (\text{A.1})$$

where $\lceil x \rceil$ gives gives the smallest integer greater than or equal to x and $\lfloor x \rfloor$ gives the largest integer less than or equal to x .

This goes until $m = T$. Therefore,

$$r_T = 2 \left(1 + \sum_{m=1}^T m \right) = 2 \left(1 + \frac{T(T + 1)}{2} \right) = T(T + 1) + 2$$

where the 1 in $\left(1 + \sum_{m=1}^T m \right)$ is accounting for the one case with $m = 0$, that is, $\{(0, T, 0, 0, 0)\}$. Note that for this proof it is not necessary to write all the possible different $\{y_0^j, n_{00}^j, n_{01}^j, n_{10}^j, n_{11}^j\}$ sets. We only wanted to count them. However, knowing (??) is going to be useful for the next proof.

A.2. Number of ‘independent’ equations with covariates: $r_{xit}(T, N_x)$

Here we proof equation (A.2), that is, that the number of different equations in the case with x_{it} covariate that takes N_x values and varies both in i and t is

$$r_{xit}(T, N_x) = 2N_x \frac{(T + N_x - 1)!}{T!(N_x - 1)!} + 2N_x \sum_{m=1}^T \sum_{q=0}^{m-1} \frac{(T - m + N_x - 1)! \left(\left\lfloor \frac{m-q}{2} \right\rfloor + N_x - 1\right)!}{(T - m)!(N_x - 1)! \left(\left\lfloor \frac{m-q}{2} \right\rfloor\right)!(N_x - 1)!} \frac{\left(\left\lfloor \frac{m-q}{2} \right\rfloor + N_x - 1\right)! (q + N_x - 1)!}{\left(\left\lfloor \frac{m-q}{2} \right\rfloor\right)!(N_x - 1)! q!(N_x - 1)!} \quad (\text{A.2})$$

It can be seen in (5.8) that now we have to count the number of different sets $\{y_0^j, x_0^j, n_{00|1}^j, \dots, n_{00|N_x}^j, n_{01|1}^j, \dots, n_{01|N_x}^j, n_{10|1}^j, \dots, n_{10|N_x}^j, n_{11|1}^j, \dots, n_{11|N_x}^j\}$ that the $j = 1, \dots, 2^{N_x(T+1)}$ possible paths can generate. $n_{01|l}^j$ is the number of $y_{t-1} = 0 \rightarrow y_t = 1$ transitions for path j given x_{it} takes the l -th value. Note that $\sum_{l=1}^{N_x} n_{00|l} = n_{00}$, so the number of 00 transitions we have for the y_t are being divided between $n_{00|1}^j, \dots,$ and $n_{00|N_x}^j$ depending on the value of x_{it} for each particular path. Therefore, we first count the number of ways n_{00} can be arranged into those N_x possible transitions without any other restriction than that (this includes that n_{00} transitions can be arranged in a way that some of the N_x new transition counters are zero). For any given value of $n_{00} = n$ this number is:

$$\frac{(n + N_x - 1)!}{n!(N_x - 1)!} \quad (\text{A.3})$$

(A.3) gives the number for a given set with $n_{00} = n$. We now have to add this for all the possible values of n_{00} . The problem and formula (A.3) are the same for n_{01} , n_{10} , and n_{11} . The number of possible sets of $\{y_0, n_{00}, n_{01}, n_{10}, n_{11}\}$ and the sets have being derived in previous appendix. There are r_T possible sets and, from equation (??), the first half of the r_T sets of $\{y_0, n_{00}, n_{01}, n_{10}, n_{11}\}$ are

$$\left\{ (0, T, 0, 0, 0), \left\{ \left\{ \left(0, T - m, \left\lfloor \frac{m-q}{2} \right\rfloor, \left\lfloor \frac{m-q}{2} \right\rfloor, q \right) \right\}_{q=0}^{m-1} \right\}_{m=1}^T \right\} \quad (\text{A.4})$$

The other half with $y_0 = 1$ can be obtained similarly, and the total number will be the number for $y_0 = 0$ multiplied by two.

Therefore, combining (A.3) and (A.4) we have that the number $r_{xit}(T, N_x)$ of possible sets of $\{y_0, x_0, n_{00|1}, \dots, n_{00|N_x}, n_{01|1}, \dots, n_{01|N_x}, n_{10|1}, \dots, n_{10|N_x}, n_{11|1}, \dots, n_{11|N_x}\}$ is given by equation (A.2) that has been written again in this appendix. The N_x comes from the number of possible values of x_0 that will give other different combinations with everything else being equal.

B. Proof of Proposition 3.1: Conditions for local identification.

Proving the local identification result in Proposition 3.1 is a direct implication of the rank of Jacobian of the system. The first section here shows that studying identification of the distribution of (P, G, H) in (3.5) is equivalent to study identification of (G, H) conditional on the first observation. Then we present several steps that simplify the system and matrices we need to analyze. This is done in order to obtain a tractable form of the Jacobian of the system. Then, we present the result and its proof about the rank of the system conditional on the first observation. Finally, using that result, we proof proposition 3.1.

B.1. Breaking the problem in two: focusing on the process conditional on the first observation.

The system of equations that defines our problem (3.5) can be expressed in terms of that system conditional on the initial observation times the distribution of the initial observation. That is $\pi = \pi_{y_0} * \Pr(y_0)$, where π_{y_0} contains the probability of each of the $\Gamma = 2^{T+1}$ paths conditional on the initial observation being y_0 . The first $\frac{\Gamma}{2}$ rows of π_{y_0} are the probability of the paths that start at $y_0 = 0$ given that and the last $\frac{\Gamma}{2}$ rows are the probability of the paths that start at $y_0 = 1$ conditional on $y_0 = 1$.¹⁶ The system is, then:

$$\pi_{y_0} = \begin{bmatrix} \pi_{y_0=0} \\ \pi_{y_0=1} \end{bmatrix} = \mathbf{A}_{y_0} \theta_{y_0} = \begin{bmatrix} \mathbf{A}_{y_0=0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{y_0=1} \end{bmatrix} \begin{bmatrix} \theta_{y_0=0} \\ \theta_{y_0=1} \end{bmatrix} \quad (\text{B.1})$$

where $\pi_{y_0=0}$ and $\pi_{y_0=1}$ are vector of dimension $\frac{\Gamma}{2} \times 1$, \mathbf{A}_{y_0} is a $\Gamma \times 2S$ matrix, θ_{y_0} is a vector of dimension $2S$, $\mathbf{A}_{y_0=0}$ and $\mathbf{A}_{y_0=1}$ are $\frac{\Gamma}{2} \times S$ matrices, $\mathbf{0}$ are $\frac{\Gamma}{2} \times S$ matrices whose elements are all zero, and $\theta_{y_0=1}$ and $\theta_{y_0=0}$ are vectors of dimension $S \times 1$. System (B.1) is simply a compact expression for two separate process: one for those observation that start with 0 and those that start with 1. The j -th element of $\pi_{y_0=0}$ and $\pi_{y_0=1}$ are respectively:

$$\pi_{j|y_0=0} = \frac{\pi_j}{\sum_{k=1}^{2^T} \pi_k}, \quad j = 1, \dots, 2^T$$

$$\pi_{j-2^T|y_0=1} = \frac{\pi_j}{\sum_{k=2^T+1}^{2^{T+1}} \pi_k}, \quad j = 2^T + 1, \dots, 2^{T+1}$$

where π_j and π_k are the elements of π in (3.5), that is, the unconditional proportions of each path. The elements of $\theta_{y_0=0}$ and $\theta_{y_0=1}$ give the probability of each type conditional on $y_0 = 0$ and $y_0 = 1$ respectively:

$$\theta_{y_0=0} = \left[\theta_{1|y_0=0}, \dots, \theta_{S-1|y_0=0}, 1 - \sum_{s=1}^{S-1} \theta_{s|y_0=0} \right]'$$

$$\theta_{y_0=1} = \left[\theta_{1|y_0=1}, \dots, \theta_{S-1|y_0=1}, 1 - \sum_{s=1}^{S-1} \theta_{s|y_0=1} \right]'$$

where

$$\theta_{s|y_0=0} = \frac{\Pr(s, y_0 = 0)}{\Pr(y_0 = 0)} = \frac{\Pr(y_0 = 0|s) * \Pr(s)}{\sum_{k=1}^S \Pr(y_0 = 0|s) * \Pr(k)} = \frac{(1 - P_s) * \theta_s}{\sum_{k=1}^S (1 - P_k) * \theta_k} \quad \text{for } s = 1, \dots, S - 1 \quad (\text{B.2})$$

$$\theta_{s|y_0=1} = \frac{P_s * \theta_s}{\sum_{k=1}^S P_k * \theta_k}, \quad \text{for } s = 1, \dots, S - 1 \quad (\text{B.3})$$

The system (B.1) contains all the information we can use to identify the distribution of (G_s, H_s) conditional on the initial observation. Once we have recovered $\theta_{y_0=0}$ and $\theta_{y_0=1}$ from that system, the distribution of P_s and the unconditional probability of each type (θ_s) can be uniquely recover from (B.2), (B.3), and the unconditional probability of the initial observation (which is a proportion

¹⁶Notice that the probability of the first $\frac{\Gamma}{2}$ paths given $y_0 = 1$ is zero because these are the paths that start at $y_0 = 0$. For the same reason the probability of the last $\frac{\Gamma}{2}$ paths given $y_0 = 0$ is zero.

of ones that we observe):

$$\Pr(y_0 = 1) = \sum_{k=1}^S P_k * \theta_k \quad (\text{B.4})$$

Once we have $\theta_{y_0=0}$ and $\theta_{y_0=1}$ and $\Pr(y_0 = 1)$, (B.2), (B.3), and (B.4) form a system of $2S - 1$ equations that uniquely identify the $2S - 1$ unknowns $(P_1, P_2, \dots, P_S, \theta_1, \dots, \theta_{S-1})$ for possible values of the parameters that are in the open interval $(0, 1)$. Furthermore, this solutions have close forms. Substituting (B.4) in (B.2) and (B.3) implies

$$(1 - \Pr(y_0 = 1)) \theta_{s|y_0=0} = \theta_s - P_s * \theta_s \quad (\text{B.5})$$

$$P_s * \theta_s = \Pr(y_0 = 1) \theta_{s|y_0=1} \quad (\text{B.6})$$

Then, substituting (B.6) in (B.5), and doing some manipulation we obtain the solution for θ_s

$$\theta_s = \theta_{s|y_0=0} * (1 - \Pr(y_0 = 1)) + \theta_{s|y_0=1} * \Pr(y_0 = 1), \text{ for } s = 1, \dots, S - 1 \quad (\text{B.7})$$

Replacing θ_s with its solution in (B.6) we obtain the solution

$$P_s = \frac{\theta_{s|y_0=1} * \Pr(y_0 = 1)}{\theta_{s|y_0=0} * (1 - \Pr(y_0 = 1)) + \theta_{s|y_0=1} * \Pr(y_0 = 1)}, \text{ for } s = 1, \dots, S - 1 \quad (\text{B.8})$$

Finally, (B.4) can be writing as $\Pr(y_0 = 1) = \sum_{k=1}^{S-1} P_k * \theta_k + P_S * \left(1 - \sum_{k=1}^{S-1} \theta_k\right)$. Substituting (B.7) and (B.8) here we can recover the solution for P_S :

$$P_S = \frac{\Pr(y_0 = 1) \left(1 - \sum_{k=1}^{S-1} \theta_{k|y_0=1}\right)}{1 - (1 - \Pr(y_0 = 1)) \sum_{k=1}^{S-1} \theta_{k|y_0=0} - \Pr(y_0 = 1) \sum_{k=1}^{S-1} \theta_{k|y_0=1}} \quad (\text{B.9})$$

This uniqueness or global invertibility in (B.2), (B.3), and (B.4) means that any non-identification problem is going to be only in (B.1). That is, if we are able to identified the distribution of (G, H) conditional on first observation, we are also able to identified the unconditional distribution of (P, G, H) .

That one-to-one map from $\left(\{\theta_{s|y_0=0}, \theta_{s|y_0=1}\}_{s=1}^{S-1}, \Pr(y_0 = 1)\right)$ to $\left(\{\theta_s\}_{s=1}^{S-1}, \{P_s\}_{s=1}^S\right)$, also shows that we can identify different values of P_s , that is, an underlying distribution of the heterogeneity in the probability of the initial observation, due to its relation with the distribution of the heterogeneity in (G, H) . If they were independent then $\theta_{s|y_0=0} = \theta_{s|y_0=1}$, and we could not identify different values of P_s but the proportions of ones we observe in the first period. That is, $\theta_{s|y_0=0} = \theta_{s|y_0=1}$ imply in (B.8) and (B.9) that $P_s = \Pr(y_0 = 1)$ for all $s = 1, \dots, S$.

B.2. Decomposition of matrix \mathbf{A}_{y_0} .

From equations (2.7) and (3.4), without the probability of the initial observation since we have conditioned on it, any element of a row j of matrices $\mathbf{A}_{y_0=0}$ and $\mathbf{A}_{y_0=1}$ is given by $G^{n_{01}^j} (1 - G)^{n_{00}^j} H^{n_{11}^j} (1 - H)^{n_{10}^j}$. From the binomial theorem we have that

$$G^{n_{01}^j} (1 - G)^{n_{00}^j} H^{n_{11}^j} (1 - H)^{n_{10}^j} = \sum_{z=0}^{n_{10}^j} \sum_{x=0}^{n_{00}^j} (-1)^x (-1)^z \binom{n_{00}^j}{x} \binom{n_{10}^j}{z} G^{(x+n_{01}^j)} H^{(z+n_{11}^j)} \quad (\text{B.10})$$

Based on this we can decompose matrix \mathbf{A}_{y_0} as the product of two matrices:

$$\mathbf{A}_{y_0} = \begin{bmatrix} \mathbf{C}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_1 \end{bmatrix} \begin{bmatrix} \mathbf{E}_1 & \dots & \mathbf{E}_S & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{E}_1 & \dots & \mathbf{E}_S \end{bmatrix} \quad (\text{B.11})$$

where $\mathbf{C} = \begin{bmatrix} \mathbf{C}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_1 \end{bmatrix}$ will contain the coefficients $\left((-1)^x(-1)^z \binom{n_{00}^j}{x} \binom{n_{10}^j}{z}\right)$ of (B.10) and \mathbf{E} will contain the corresponding G , H and P terms. The matrix \mathbf{C} does not depend on the value of the parameters and, therefore, it will be unique for a given T .

\mathbf{E}_s is the following vector

$$\mathbf{E}'_s = \begin{bmatrix} 1 & G_s & \dots & G_s^T & H_s & G_s H_s & \dots & G_s^{T-1} H_s & H_s^2 & \dots & G_s^{T-2} H_s^2 & \dots & H_s^{T-1} & G_s H_s^{T-1} & H_s^T \end{bmatrix} \quad (\text{B.12})$$

of dimension

$$e_T = \frac{(T+1)(T+2)}{2} \quad (\text{B.13})$$

Notice that e_T is the triangular number $(T+1)$. For instance, with $T=2$

$$\mathbf{E}_s = \begin{bmatrix} 1 & G_s & G_s^2 & H_s & G_s H_s & H_s^2 \end{bmatrix}'$$

Define \mathbf{C}_0 as $\frac{\Gamma}{2} \times e_T$ matrix whose row j have the binomial coefficients from the path (the binary number with $T+1$ digits) that correspond with the decimal number $(j-1) : j = 1, \dots, \frac{\Gamma}{2}$. For instance, the third row with $T=2$ corresponds with the path 010, which is the three-digit binary number that represents the decimal number 2. This way of using the corresponding decimal numbers to order the paths and rows of \mathbf{C}_0 , also implies the order of the elements of vector \mathbf{E}_s . Each row j in \mathbf{C}_0 contains the coefficients of the different terms of (B.10) plus the zeros needed to filling the rest of the cells for those elements in \mathbf{E}_s that do not appear in the probability of path j . A coefficient $\left((-1)^x(-1)^z \binom{n_{00}^j}{x} \binom{n_{10}^j}{z}\right)$ is completely defined by j , x and z , and it is in row j and column

$$(Z + n_{11}^j)(T+2) - \frac{(z + n_{11}^j)(z + n_{11}^j + 1)}{2} + x + 1 + n_{01}^j \quad (\text{B.14})$$

of matrix \mathbf{C}_0 .

Define \mathbf{C}_1 the same way as \mathbf{C}_0 , but $j = \frac{\Gamma}{2} + 1, \dots, T$. Each coefficient of (B.10) is in column given by (B.14) and row $j - \frac{\Gamma}{2}$. Then,

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_0 & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_1 \end{bmatrix} \quad (\text{B.15})$$

The dimension of \mathbf{C} is $\Gamma \times 2e_T$ and the dimension of each sub-matrix \mathbf{C}_0 and \mathbf{C}_1 is $\frac{\Gamma}{2} \times e_T$. From

(B.10) and (B.14) matrix \mathbf{C} can be easily computed for any given T . For example, with $T = 2$

$$\mathbf{C} = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 & 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \quad (\text{B.16})$$

with dimension 8×12 .

B.3. Eliminating redundancies in \mathbf{A}_{y_0} .

As stated in Proposition 2.2 only $r_T (= T(T+1)+2)$ of the $\Gamma (= 2^{T+1})$ possible paths are distinct paths. Therefore, \mathbf{A}_{y_0} can not have a rank bigger than r_T since $\Gamma - r_T$ rows in \mathbf{A}_{y_0} are repetitions of rows whose paths are the same. Here eliminate those redundancies in \mathbf{A}_{y_0} since it is the rank of it what will define the rank of the system and the rank of its Jacobian will be the rank of J . Let us denote the matrix without redundancies with the subscript r .

First we take from \mathbf{C} those rows j that correspond to a path j that is not different from a previous path. This means that the number of rows of \mathbf{C}_r is r_T . Secondly we reduce the number of columns in \mathbf{C} that are zero or can be expressed as linear combinations of other columns. This means that we are eliminating $2T$ columns (T in each sub-matrix \mathbf{C}_0 and \mathbf{C}_1) so that the number of columns of \mathbf{C}_r equals r_T too. This column reduction requires the corresponding adjustment in E_s .

In \mathbf{C}_0 this only requires to eliminate the T columns that are zero. These columns correspond with the T elements in E_s that are only a function of H_s . That is, we eliminate H_s, H_s^2, \dots , and H_s^T from E_s when using it in the part of the system that gives the probability for those paths starting with 0. These elements are part of (B.10) for the paths that start at $y_0 = 1$ but and H_s can not be alone an element of (B.10) when $y_0 = 0$. Then,

$$\mathbf{E}'_{s,0r} = \begin{bmatrix} 1 & G_s & \dots & G_s^T & G_s H_s & \dots & G_s^{T-1} H_s & G_s H_s^2 & \dots & G_s^{T-2} H_s^2 & \dots & G_s H_s^{T-1} \end{bmatrix} \quad (\text{B.17})$$

In \mathbf{C}_1 , in addition to a column that is zero and corresponds to element G_s^T in vector E_s , there are $T - 1$ columns that are linear combinations of other columns. These $T - 1$ columns to be eliminated from \mathbf{C}_1 correspond to $\{G_s^{T-i} H^i\}_{i=1}^{T-1}$ in E_s . We eliminate $\{G_s^{T-i} H^i\}_{i=1}^{T-1}$ from E_s and replace $\left\{ \left\{ G_s^{T-i} H_s^j \right\}_{j=0}^{i-1} \right\}_{i=1}^{T-1}$ by $\left\{ \left\{ G_s^{T-i} H_s^j (1 - H_s^{i-j}) \right\}_{j=0}^{i-1} \right\}_{i=1}^{T-1}$. This reflects the fact that, for paths starting at $y_0 = 1$, G_s^T can not be part of (B.10), and G_s with any exponent will only appear in (B.10) if there is at least a $(1 - H)$, given that G is $\Pr(y_t = 1 | y_{t-1} = 0)$. Thus, the vector $\mathbf{E}'_{s,1r}$ is of dimension $\frac{r_T}{2}$, its typical element is $G_s^{T-i} H_s^j (1 - H_s^{i-j})$ for $i = 1, \dots, T - 1$ and $j = 0, \dots, i - 1$, which is in position $T + i + 1 + \mathbf{1}\{j > 0\} \sum_{k=0}^{j-1} (T - k)$ in the vector. That is,

$$\mathbf{E}'_{s,1r} = \begin{bmatrix} 1 & G_s (1 - H_s^{T-1}) & \dots & G_s^{T-1} (1 - H_s) & H_s & G_s H_s (1 - H_s^{T-2}) & \dots & G_s^{T-2} H_s (1 - H_s) & H_s^2 \\ G_s H_s^2 (1 - H_s^{T-3}) & \dots & G_s^{T-3} H_s^2 (1 - H_s) & \dots & \dots & H_s^{T-2} & G_s H_s^{T-2} (1 - H_s) & H_s^{T-1} & H_s^T \end{bmatrix} \quad (\text{B.18})$$

For example, with $T = 2$ and $S = 2$:

$$\begin{aligned}
& \begin{bmatrix} \pi_{1|y_0=0} \\ \pi_{2|y_0=0} \\ \pi_{3|y_0=0} \\ \pi_{4|y_0=0} \\ \pi_{5|y_0=1} \\ \pi_{6|y_0=1} \\ \pi_{7|y_0=1} \\ \pi_{8|y_0=1} \end{bmatrix} = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 & 0 & -1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} * \\
& \begin{bmatrix} 1 & 1 & 0 & 0 \\ G_1 & G_2 & 0 & 0 \\ G_1^2 & G_2^2 & 0 & 0 \\ H_1 & H_2 & 0 & 0 \\ G_1H_1 & G_2H_2 & 0 & 0 \\ H_1^2 & H_2^2 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & G_1 & G_2 \\ 0 & 0 & G_1^2 & G_2^2 \\ 0 & 0 & H_1 & H_2 \\ 0 & 0 & G_1H_1 & G_2H_2 \\ 0 & 0 & H_1^2 & H_2^2 \end{bmatrix} \begin{bmatrix} \theta_{1|y_0=0} \\ 1 - \theta_{1|y_0=0} \\ \theta_{1|y_0=1} \\ 1 - \theta_{1|y_0=1} \end{bmatrix} \\
& = \begin{bmatrix} 1 & -2 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 & 0 \\ G_1 & G_2 & 0 & 0 \\ G_1^2 & G_2^2 & 0 & 0 \\ G_1H_1 & G_2H_2 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & G_1(1-H_1) & G_2(1-H_2) \\ 0 & 0 & H_1 & H_2 \\ 0 & 0 & H_1^2 & H_2^2 \end{bmatrix} \begin{bmatrix} \theta_{1|y_0=0} \\ 1 - \theta_{1|y_0=0} \\ \theta_{1|y_0=1} \\ 1 - \theta_{1|y_0=1} \end{bmatrix} \\
& \tag{B.19}
\end{aligned}$$

where the three matrices in the last line are respectively denoted by $\mathbf{C}_r (= \begin{bmatrix} \mathbf{C}_{0r} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{1r} \end{bmatrix})$, \mathbf{E}_{y_0r}

($= \begin{bmatrix} \mathbf{E}_{1,0r} & \mathbf{E}_{2,0r} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{E}_{1,1r} & \mathbf{E}_{2,1r} \end{bmatrix}$), and θ_{y_0} .

B.4. Isolating the unknown parameters

It is important to note that \mathbf{C} and \mathbf{C}_r do not depend on unknown parameters. We can construct \mathbf{C} and \mathbf{C}_r and calculate the rank of \mathbf{C} for any given T , using (B.10), (B.14) and indications in previous subsection B.3. Obviously, the rank of \mathbf{C} is equal to the rank of \mathbf{C}_r . Table 2.2 reports $rank(\mathbf{C})$, for $T = 2, \dots, 23$. For all those values of T , the rank of \mathbf{C} is the number of equations that are different in the system, r_T :

$$r_T = T(T + 1) + 2 = rank(\mathbf{C}) = rank(\mathbf{C}_r) \tag{B.20}$$

r_T is also the dimension of the square matrix \mathbf{C}_r . That is, \mathbf{C}_r is a matrix of full rank and we

can invert it. Then, our system conditional in the first observation:

$$\pi_{y_0 r} = \mathbf{C}_r \mathbf{E}_{y_0 r} \theta_{y_0}$$

is equivalent to

$$\mathbf{C}_r^{-1} \pi_{y_0 r} = \mathbf{E}_{y_0 r} \theta_{y_0} \quad (\text{B.21})$$

where $\mathbf{C}_r^{-1} = \begin{bmatrix} \mathbf{C}_{0r}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{1r}^{-1} \end{bmatrix}$ and $\pi_{y_0 r}$ are the proportions, conditional on y_0 , for the r_T paths that are different. The advantage of (B.21) is that the right hand side contains only unknown values. The two subsystems in (B.21), one conditional on $y_0 = 0$ and the other conditional on $y_0 = 1$, are

$$\mathbf{C}_{0r}^{-1} \pi_{y_0=0r} = \begin{bmatrix} 1 & \dots & 1 \\ G_1 & \dots & G_S \\ \vdots & \dots & \vdots \\ G_1^{T-i} H_1^j & \dots & G_S^{T-i} H_S^j \\ \vdots & \dots & \vdots \\ G_1 H_1^{T-1} & \dots & G_S H_S^{T-1} \end{bmatrix} \begin{bmatrix} \theta_{1|y_0=0} \\ \vdots \\ \theta_{S-1|y_0=0} \\ 1 - \sum_{s=1}^{S-1} \theta_{s|y_0=0} \end{bmatrix} \quad (\text{B.22})$$

$$\mathbf{C}_{1r}^{-1} \pi_{y_0=1r} = \begin{bmatrix} 1 & \dots & 1 \\ G_1 (1 - H_1^{T-1}) & \dots & G_S (1 - H_S^{T-1}) \\ \vdots & \dots & \vdots \\ G_1^{T-i} H_1^j (1 - H_1^{i-j}) & \dots & G_S^{T-i} H_S^j (1 - H_S^{i-j}) \\ \vdots & \dots & \vdots \\ H_1^T & \dots & H_S^T \end{bmatrix} \begin{bmatrix} \theta_{1|y_0=1} \\ \vdots \\ \theta_{S-1|y_0=1} \\ 1 - \sum_{s=1}^{S-1} \theta_{s|y_0=1} \end{bmatrix} \quad (\text{B.23})$$

It is clear from a direct inspection of (B.22) and (B.23) that the first equation in each of these two systems is trivially satisfied for any value of the parameters since it is the sum of the probability of point of support s which by definition is always equal to one. Therefore, although (B.21) contains r_T different equations, only $r_T - 2$ equations restrict the value of the unknowns. This correspond to the fact that elements in $\pi_{y_0=0}$ and $\pi_{y_0=1}$ sum one, so one of them can be expressed as a linear combination of all the other elements inside each subsystem. Given this, in what follow when we refer to the these systems (B.21), (B.22), and (B.23), and to matrix $\mathbf{E}_{y_0 r}$ we are referring to their formulation without the first elements that trivially sum one. Then the dimension of $\mathbf{E}_{y_0 r}$ is $(r_T - 2) \times S$.

B.5. The rank of \mathbf{J}_r matrix

The identification result that we try to proof is based on \mathbf{J} (the Jacobian matrix of system (3.5)) having rank greater or equal than the number of unknowns. This is equivalent to the Jacobian of (B.21) having rank greater or equal than the number of unknowns in this system. We denote the latter by \mathbf{J}_r . \mathbf{J}_r is a matrix of dimension $(r_T - 2) \times (4S - 2)$ composed by the following parts:

- First S columns that contain the derivatives with respect to G_1, \dots, G_S , whose general form is

$$\begin{bmatrix} \theta_{s|y_0=0} \\ \cdot \\ (T-i)G_s^{T-i-1}H_s^j\theta_{s|y_0=0} \\ \cdot \\ H_s^{T-1}\theta_{s|y_0=0} \\ (1-H_s^{T-1})\theta_{s|y_0=1} \\ \cdot \\ (T-i)G_s^{T-i-1}H_s^j(1-H_1^{i-j})\theta_{s|y_0=1} \\ \cdot \\ 0 \end{bmatrix} \quad (\text{B.24})$$

- Next S columns that contain the derivatives with respect to H_1, \dots, H_S , whose general form is

$$\begin{bmatrix} 0 \\ \cdot \\ jG_s^{T-i}H_s^{j-1}\theta_{s|y_0=0} \\ \cdot \\ (T-1)G_sH_s^{T-2}\theta_{s|y_0=0} \\ -(T-1)G_sH_s^{T-2}\theta_{s|y_0=1} \\ \cdot \\ G_s^{T-i}(jH_s^{j-1}(1-H_s^{i-j})-(i-j)H_s^{i-1})\theta_{s|y_0=1} \\ \cdot \\ TH_s^{T-1}\theta_{s|y_0=1} \end{bmatrix} \quad (\text{B.25})$$

- Last $2(S-1)$ columns that contain the derivatives with respect to $\theta_{1|y_0=0}, \dots, \theta_{S-1|y_0=0}, \theta_{1|y_0=1}, \dots, \theta_{S-1|y_0=1}$, whose general form is

$$\begin{bmatrix} G_s - G_S \\ \cdot \\ G_s^{T-i}H_s^j - G_S^{T-i}H_S^j \\ \cdot \\ G_sH_s^{T-1} - G_SH_S^{T-1} \\ G_1(1-H_1^{T-1}) - G_S(1-H_S^{T-1}) \\ \cdot \\ G_s^{T-i}H_s^j(1-H_s^{i-j}) - G_S^{T-i}H_S^j(1-H_S^{i-j}) \\ \cdot \\ H_s^T - H_S^T \end{bmatrix} \quad (\text{B.26})$$

For example, with $T = 2$ and $S = 2$:

$$\mathbf{J}_r = \begin{bmatrix} \theta_{1|y_0=0} & (1 - \theta_{1|y_0=0}) & 0 & 0 & G_1 - G_2 & 0 \\ 2G_1\theta_{1|y_0=0} & 2G_2(1 - \theta_{1|y_0=0}) & 0 & 0 & G_1^2 - G_2^2 & 0 \\ H_1\theta_{1|y_0=0} & H_2(1 - \theta_{1|y_0=0}) & G_1\theta_{1|y_0=0} & G_2(1 - \theta_{1|y_0=0}) & G_1H_1 - G_2H_2 & 0 \\ (1 - H_1)\theta_{1|y_0=1} & (1 - H_2)(1 - \theta_{1|y_0=1}) & -G_1\theta_{1|y_0=1} & -G_2(1 - \theta_{1|y_0=1}) & 0 & * \\ 0 & 0 & \theta_{1|y_0=1} & (1 - \theta_{1|y_0=1}) & 0 & H_1 - H_2 \\ 0 & 0 & 2H_1\theta_{1|y_0=1} & 2H_2(1 - \theta_{1|y_0=1}) & 0 & H_1^2 - H_2^2 \end{bmatrix} \quad (\text{B.27})$$

$$* = G_1(1 - H_1) - G_2(1 - H_2) \quad (\text{B.28})$$

Since we are trying to derive the minimum number of periods needed for identifying a distributions with S points of support, we first look at the case where S is not limiting the rank of the matrix. Therefore, we consider here a case where \mathbf{J}_r is a squared matrix: $4S - 2 = r_T - 2$. If squared \mathbf{J}_r matrix has full rank, this will give the identification condition.

\mathbf{J}_r depends on the value of unknown parameters, and so does the determinant of it $\det(\mathbf{J}_r)$. Therefore, by simply looking at its general form we can not conclude whether $\det(\mathbf{J}_r)$ is different from zero for all the possible values of the parameters. However, it is not difficult to see that if we evaluate $\det(\mathbf{J}_r)$ at values of the parameters where there is no special relations between the different parameters and points of support all the rows and columns in \mathbf{J}_r are linearly independent and, therefore $\det(\mathbf{J}_r) \neq 0$ when evaluated at those values. Furthermore, simulating many times the matrix \mathbf{J}_r with random draws for the the P_s 's, G_s 's and H_s 's we found for all those values that squared \mathbf{J}_r has $\det(\mathbf{J}_r) \neq 0$ that is, full rank and, therefore the rank of \mathbf{J}_r is given by: $r_T - 2$. Of course this only shows that \mathbf{J}_r has full rank for those particular numbers tried on the simulations. However finding even only one point for which $\det(\mathbf{J}_r) \neq 0$ is going to be crucial to prove a result about the rank of \mathbf{J}_r in general. The argument is as follows.

Firstly, it is important to note that the equations in system (B.21) are polynomial functions and, therefore $\det(\mathbf{J}_r)$ is also a polynomial function $\mathbf{R}^{4S-2} \rightarrow \mathbf{R}$. A polynomial function is either identically zero; that is, it is zero for all values at which that function is evaluated, or the set of values at which is zero, (its roots) is of measure zero in \mathbf{R}^{4S-2} . This result is proved in Lemma 1.1 of Eisenfeld (1986). According to this result, if the polynomial function is not identically zero, then it is different from zero almost everywhere. Therefore, using Lemma 1.1 of Eisenfeld (1986), it is enough to have found a value of the parameters such that $\det(\mathbf{J}_r) \neq 0$ at that particular point to conclude that the $\det(\mathbf{J}_r) \neq 0$ almost everywhere. Putting it in different words, given that there are points at which \mathbf{J}_r has full rank, the set of values of the unknown parameters for which squared \mathbf{J}_r does not have full rank is a set of measure zero.

Given this result, we can conclude that for any given S ,

$$\text{rank}(\mathbf{J}_r) = \min(r_T - 2, 4S - 2)$$

almost everywhere.

B.6. Proof of proposition 3.1

That (3.9) in proposition 3.1 is a sufficient condition for identification is a direct application of the general inverse function theorem and the result about the rank of the Jacobian we have shown above. For local point identification of $(G, H)|_{y_0}$ the inverse function theorem requires that the rank of \mathbf{J}_r be equal to the number of unknown parameters. As shown in B.5, the rank of \mathbf{J}_r is equal to $\min(r_T - 2, \text{number of unknown parameters of the distribution conditional on the first$

observation). Therefore, the requirement for this case is that the number of unknowns be smaller than or equal to $r_T - 2$, that is $4S - 2 \leq T(T + 1)$. From here we obtain the sufficient condition to identify a distribution of $(G, H) | y_0$ in equation (3.9). To proof that this condition is also sufficient for identification of the distribution of (P, G, H) with S points of support, it is enough to recall that $(G, H) | y_0$ and the observed probability of the initial observation define a unique value for the parameters of the unconditional distribution of (P, G, H) , as shown in B.1.

To proof the necessity of condition (3.9) we use Theorem 5.A.1. in Appendix to Chapter 5 in Fisher (1966). That Theorem states that having the rank being equal to the number of unknowns is a necessary condition for a local identification of a solution if that solution is a regular point. A point is defined as regular when for all points in a sufficiently small neighborhood of it the Jacobian has the same rank as in the point (see definition 5.A.1 in Appendix to Chapter 5 in Fisher, 1966). As shown in B.5 the rank of the Jacobian is constant for all points for which (3.9) is a sufficient condition (that is, all points except a set of point of zero measure). Therefore, for those points it also a necessary condition for identification.