# Life History Analysis in Demography: Implications for Teaching and Research

**Fernando Rajulton**
Department of Sociology
The University of Western Ontario
London, Ontario, Canada

*Abstract*
In demography, a retrospective observation plan has become a common data collection procedure, and almost all surveys done — in recent times collect life history data. To justify the collection of life history information, all three of its aspects — namely, the timing, sequence and number — should be considered in a meaningful analysis. A multivariate analysis of the timing of one or two events using a number of covariates without consideration of the sequence or number of events is not, and cannot be, a life history analysis. By bringing out relevant points regarding the fundamental assumptions in life history analyses, this paper aims at contributing toward developing theories of change and procedures of estimation and testing. In particular, three stochastic models — Markov, semi-Markov and non-Markov — are discussed in detail. Other possible models, including those of diffusion as well as of unobserved heterogeneity, are suggested.

*Résumé*
En démographie, le plan d'observation rétrospectif est devenu un mode habituel de collecte de données et presque toutes les enquêtes effectuées récemment recueillent les données du cycle de vie. Pour que cette démarche se justifie, les trois aspects du cycle — moment, séquence et nombre — devraient être pris en considération dans toute analyse significative. L'analyse multifactorielle du moment de survenue d'un ou de deux événements qui utilise un certain nombre de covariables, mais sans s'intéresser à la séquence ou au nombre d'événements, ne mérite pas le nom d'analyse du cycle de vie. Le présent article vise à contribuer à l'élaboration des théories de changement et aux modes d'évaluation et de vérification en faisant ressortir les éléments pertinents des hypothèses fondamentales des analyses de ce type. Trois processus stochastiques particuliers sont décrits en détails: la chaîne de Markov, les modèles semi-markovien et non markovien. Les modèles de diffusion et d'hétérogénéité non observée figurent également parmi les autres modèles suggérés.

*Key Words*:  life history data, stochastic processes, unobserved heterogeneity, parametric models

## Introduction

Life is marked by a sequence of events. Individuals are born, enter school, graduate, are employed, marry, give birth, migrate, are divorced, become widowed, and ultimately die. The occurrences of these events and their consequences on populations and societies have been the focus of research in the social sciences. Because events are qualitative changes that occur at specific points in time, and because an individual's life can be characterized by a particular sequence of events, the best way to study events (as well as their causes and effects) is through complete (or partially complete) information on the number, timing and sequence of events. Data which provide such information are called event history or

life history data. This is in contrast with other incomplete types, such as panel data, which record events at a set of arbitrary times — usually equal-spaced and obtained through censuses or surveys at regular intervals — and event-count data, which record the number of events in an interval. Event histories provide the best information (depending on how complete the enquiry is) on each individual's sample path, which traces the course of events and the different states an individual occupies along with the duration spent in each state. Thus, properly speaking, a set of event history data consists of all sample paths of all individuals in a sample.

Event history data are collected through a prospective or retrospective observation plan (Pressat, 1969). In a prospective plan, members of a cohort are observed and their demographic experience is recorded as soon as it occurs. A few countries keep population registers wherein a record of each individual is maintained and updated when events occur. This is an ideal situation that provides a complete continuous time recording of all events. In a retrospective plan, individuals are asked at certain times through a survey about their experiences in the past. Its major difference from a prospective observation plan is that experiences of only a subcohort can be recorded through this plan; some individuals are necessarily excluded from recording because of death or emigration.

There are two approaches in measuring an event through a retrospective plan. First, individuals are asked about the events, either all or selected ones, experienced from a specific time in the past up to the time of interview or to an earlier time. Second, certain events (such as migration) are measured indirectly by comparing the current status (as current region of residence) with the one at some previous point in time (such as region of birth or residence one year before). In migration analysis, these two approaches yield two different sets of data: the first records events (migrations), while the second records individuals who experience events (migrants). The distinction between these two is unnecessary in analyses of other events, where the number of events and the number of individuals who experience them are identical. The present paper is confined to the first type of data — records of events.

Even though retrospective data have many advantages, they also have a few typical drawbacks. Two of the serious ones are: (*a*) a sample may not be representative of all persons in a particular cohort of interest because of the death or emigration of some members before the time of enquiry; the individuals thus excluded from observation may be a select group; and (*b*) recall error or systematic mistiming of events can yield spurious trends; it has been established, however, that though errors in report of timing are

2

quite common, the logical sequence of events is usually reported correctly (Tuma and Hannan, 1984).

In demography, a retrospective observation plan has become a common data collection procedure, and almost all surveys done in recent times collect life history data. For example, the World Fertility Surveys in the 1970's and the National Migration Surveys conducted in a few developing countries in the ESCAP region in the 1980's collected life history data on specific demographic phenomena such as migration, nuptiality and fertility. In spite of this, however, analyses are still being carried out as if panel observation plans had been used! At best, multivariate techniques are used to analyze a single event or two selected from a sequence of different events, thus neglecting to consider the course of events in an individual's life. To justify the collection of life history information (no doubt done at a very high cost!), all its three aspects — namely, the timing, sequence and number — should be considered in a meaningful analysis. Contrary to claims made by some research papers appearing in refereed journals, a multivariate analysis of timing of one or two events using a number of covariates but without consideration of the sequence or number of events is not, and cannot be, a life history analysis. This present state of affairs can be improved by developing (*a*) theories of change and (*b*) procedures of estimation and testing. This paper aims at contributing toward this development by bringing out relevant points regarding two fundamental assumptions in life history analyses: (*a*) that a specific stochastic process generates events, which can be appropriately analyzed; and, (*b*) that certain characteristics of individuals, as well as of contexts, affect change processes.

*Stochastic Models of Life History Analysis*
The first assumption, that a specific stochastic process generates events, touches on the basic modelling procedure: whether a model should be deterministic or stochastic; that is, whether the effect of any change in a system can be predicted with certainty or not. It is common knowledge that no social system is fully determined, controllable or predictable, and that no human behaviour warrants deterministic predictions. Nobody can predict with certainty, for example, whether a son will achieve his father's status, when a man/woman will decide to change his/her job or marital status, whether a chance encounter between members of a group will lead to a diffusion of new social trends, and so on. This uncertainty implicit in social systems can be taken into account only by introducing probability distributions into a model. Simply stated, the equations of a model should include random variables. A model which accounts for a large element of chance in the process under study and which contains variables that cannot

be fully controlled or determined prior to observation (and hence have a probability other than one) is called a stochastic model (Neyman called it "dynamic indeterminism"). Predictions, then, become probabilistic since probabilities are assigned to various possible future statuses of the process. Probabilistic predictions are important especially for normative behaviour with which demographers generally deal and which is mostly associated with the concept of risk.

Apart from the question of uncertainty, there are many other reasons, both substantive and technical, for using stochastic models (for details, see Tuma and Hannan, 1984). To cite a few, stochastic models: (1) enable us to extend the theory of rational decisions to those situations where the outcomes of decisions or the circumstances influencing the outcomes are not known with certainty; (2) enhance our understanding of systemic relationships beyond simple "chains of causality" (see below for a discussion on this topic); (3) provide opportunities for including unobservables or unmeasurables in analysing the influence of covariates (discussed below); (4) can explain the evolution of a distribution, even when the initial distribution is uniform; this is in contrast to a deterministic model, which can explain a change from some initial distribution but cannot explain how the initial distribution itself arose in the first place; (5) are a better approach to the dynamics of variability in distributions, since many social processes are often revealed more clearly in variances (of behaviours) than in averages; for example, the shape of an income distribution rather than its mean; or, variances in timing of progression into higher parities rather than total fertility rate, and so on; (6) make possible a joint study of qualitative and quantitative changes; for example, interdependent linkages between income and changes in marital status, or between type of city and crime rates.

Basic stochastic processes or probabilistic laws governing the occurrence of events can be inferred from the observed distributions. For example, to say that the number of events at time $t$ depends on the number of events that already occurred implies a contagion or diffusion process. Or to say that the number of events falls into different categories, each with a specific process but with different parameters, points to a problem of heterogeneity. And to say that the parameters themselves are functions of time indicates a process which is non-stationary.

A stochastic process is simply a collection of random variables that describes the evolution of a process over time. We shall denote by $X(t)$ the number of events at time $t$, or, in general, the number of individuals who have experienced a given event at time $t$. If a process happens to be

4

bidimensional (for example, the number married at any time depends on both the formation and dissolution of marriages), we shall denote the corresponding number by a vector $X(t) = \{X_1(t), X_2(t)\}$, which represents two related processes. Extension to the multidimensional case is straightforward.

In order to deal with life history data, the above concept can be generalized through the following: (1) Instead of associating $X(t)$ with the number of events at time $t$, we can think of life history data as describing the values of a qualitative variable $X(t)$ within some observation period $T$. The set of distinct values of $X(t)$ can be denoted by a state space $S$, and the number of distinct (mutually exclusive) values are the size of the space. For example, in a marital status analysis, the state space $S$ may consist of five marital states and their associated values as follows: never married = 1, cohabiting = 2, presently married = 3, widowed = 4, and divorced = 5. (2) It is also possible to include absorbing states, which once entered, cannot be left, depending on the type of analysis. In a study of parity progression, for example, an analyst may decide to introduce an absorbing state "sterilization". (3) In addition, let $t$ denote some continuous time parameter, say $0 < t < \omega$ where $\omega$ is the last time point considered for analysis; in most demographic studies, the time parameter often represents age or duration (or seniority) since an event-origin. In order to consider the sequence of events of transitions from one state to another, we shall attach a subscript $n$ to the time parameter (that is, $t_n$) denoting the timing of the $n$-th event. With these definitions, events n denote changes in $X(t)$, the random variable $t_n$ is the time of the $n$-th event, and the random variable $X(t_n)$ refers to the $n$-th state occupied.

What is important to note is that in many social processes, $X(t)$ for different $t$ can be mutually dependent. Therefore, relationships among them assume a prominent role in stochastic analytic procedures. In fact, stochastic processes are broadly described according to the nature of dependence among the random variables. The basic concepts of a few processes relevant to life history analysis are spelt out here. Mathematical expositions are avoided as much as possible; readers who are interested in mathematical detail should consult texts on stochastic processes.

A Markov Model

The most commonly used stochastic process to study the timing and sequence of events is the Markov process. If $X(t)$, $t \, \varepsilon \, T$ is a stochastic process such that given the value of $X(s)$, the values of $X(t)$, $t > s$, do not depend on the values of $X(u)$, $u < s$, then the process is Markovian. This implies a simple dependency of events — the occurrence of the event of

interest (or, transition from one state to another) depends directly on that of the preceding event (or state), and only on it. Such a relationship gives rise to a conditional probability of $X(t)$ given $X(s)$, the manner in which $X(s)$ was reached being of no consequence.

In order to simplify notations, let $X(t)$ denote a state $k$ and $X(s)$ another state $j$. Then, the implied conditional probability can be simply denoted by $P_{jk}$ $(s,t)$, which is interpreted as the probability that an individual will be found in state $k$ at time $t$, given that the individual occupies state $j$ at time $s$. This probability is called a state transition probability. When the state space is discrete, the process is called a Markov Chain. If one considers the state space to consist of the five marital states mentioned above, then $P_{35}(29,38)$ would denote the probability that an individual who is presently married at age 29 will be divorced at age 38.

In the case of Markov chains, it would be useful to note the following :

1. In the above definition, there is no probability attached to a single state $j$ or $k$, but to a pair of states $(j,k)$ at two successive time points.

2. If the probability is dependent only on the time difference $(t-s)$, then the process is said to be homogeneous with respect to time; and if the probability is dependent on either $s$ or $t$, then the process is non-homogeneous with respect to time.

3. Two successive events lead to a one-step transition. Since life history data contain sequence of events, these one-step transitions can usually be obtained. We can generalize the concept to multi-step transitions, say state $j$ as the $n$-th event and state $k$ as the $(n+m)$-th event. This involves an $m$-step transition, and its probability is denoted by $P_{jk}^{(m)}(s,t)$. However, since it is a Markov process, each event is dependent only on the immediately previous event, which leads to the so-called Chapman Kolmogorov equation.

4. Since $P_{jk}(s,t)$ is a probability, it should satisfy the following conditions. $P_{jk}(s,t) > 0$ and $\Sigma P_{jk}$ $(s,t) = 1$ for all $j$, and summing over $k$. This property is called a stochastic property, and a matrix containing such probabilities with origin states as rows and destination states as columns is called a stochastic matrix. This matrix has non-negative elements and unit row sums. The following is an example of a stochastic matrix containing

6

probabilities of transitions between ages 29 and 30 in a Markov system of marital states described earlier in the text (that is, never married = 1, cohabiting = 2, presently married = 3, widowed = 4, and divorced = 5, with no absorbing state):

Destination

| Origin | 1 | 2 | 3 | 4 | 5 |
|--------|------|------|------|------|------|
| 1 | .324 | .654 | .001 | .017 | .004 |
| 2 | 0 | .960 | .003 | .035 | .002 |
| 3 | 0 | .406 | .483 | .010 | .101 |
| 4 | 0 | .661 | .001 | .330 | .008 |
| 5 | 0 | .322 | .276 | .002 | .400 |

The diagonal elements are retention probabilities while the off-diagonal elements are transition probabilities from the origin state at age 29 to the destination state before age 30.

5. Finally, a technical point: A Markov model is such that the transition rates (or instantaneous rates) are constant across time difference *(t-s)* and across individuals. Mathematically speaking, it is this constancy of rates which gives rise to the model's simplicity. The constancy of rates may be easily admitted when the interval *(t-s)* is small, but becomes increasingly inadmissible for longer intervals. In fact, the constancy of rates implies an exponential distribution of timing of transitions.

The Markov model has been applied to a wide variety of social phenomena as in labour mobility (Blumen *et al.,* 1955), change in attitudes (Coleman, 1964), and collective violence (Spilerman, 1970) and, in more recent times, to the development of multistate demography (Schoen and Land, 1979; Willekens and Rogers, 1978; Hoem and Funck Jensen, 1982). However, a simple model rarely describes reality well, and doubts have been raised as to whether any social process obeys the Markovian assumption at all. At the same time, however, it has been repeatedly pointed out that the final results of a cumbersome analysis of social processes using more complex frameworks are not much different from those obtained with the Markovian framework (Bartholomew, 1967). In the light of these opposing views, what Howard (1971;4) says makes sense; that is,

No experiment can ever show the ultimate validity of the Markovian assumption; hence no physical systems can ever be classified absolutely as either Markovian or non-Markovian — the

important question is whether the Markov model is useful. If the Markovian assumption can be justified, then the investigator can enjoy analytical and computational convenience not often found in complex models.

## A Semi-Markov Model

In line with arguments calling for more realism in life history analysis, various extensions of the Markov model have been tried by including, for example, population heterogeneity, time-dependence (duration, seniority or experience), and so on. Of particular interest in analyzing life history data is the influence of duration spent in each state before making a transition. When the time parameter $t$ denotes duration other than age, it becomes meaningless to hold on to the Markov model because the distributions of waiting times in states of a Markov process are exponential; in other words, these distributions are without memory and, therefore, not duration-dependent (see Ross, 1983). If duration is an important variable in an analysis, as very often it is, then the Markovian framework has to be reformulated. Following Feller's ideas (1950, 1966), Mode (1982, 1985) suggested ways of defining the probabilities involved in duration analysis directly on sample paths. Mode's approach is preferred to other approaches (Ginsberg, 1971; Hoem, 1972), as it easily leads to computer algorithms for calculating the basic probabilities involved.

A semi-Markov model considers changes in states according to a Markov chain, but allows the distribution of time intervals between successive transitions to be arbitrary (arbitrary in the sense that the distribution can be other than exponential; if it happens to be exponential, then the semi-Markov model becomes indentical to the Markov model) and to depend on the state of origin (as in a Makov chain) as well as on the state of destination (unlike a Markov chain). Thus, the dependency of events in a semi-Markovian framework is a modified form of a Markovian one: a transition from one state to another depends both on the origin state and the destination state, and on the length of duration in the origin state.

To express the above concept in stochastic terminology, we need to consider the timing and sequence of events as follows. Let $s_i$, $i > 0$ be the state in $S$ in an individual's sample path; and let $y_i$ be the time taken to go from $s_{i-1}$ to $s_1$, $i > 1$; in other words, $y_i$ is the sojourn time (duration) in state $s_{i-1}$. If we denote by $w_i = (s_i, y_i)$ the $i$-th transition among the states of $S$, $i > 1$, then $w = (w_0, w_1, \ldots)$ denotes an individual's sample path consisting of sequences of transitions. Note that $w_i$ is defined for $i > 1$. Since age is an important variable in demographic analysis, we can also introduce $w_0 = (s_0, x_0)$ to denote the pair describing the $0$-th step, where

$x_0$ denotes the age at which the state was entered; $w_0$ would simply mean that an individual aged $x_0$ enters state $s_0$, and then $w_1 = (s_1, y_1, x_0)$ follows as before. We shall, however, focus our attention in this paper on sample paths without considering the age at entry into each state. Interested readers can consult Mode (1985) for details.

Let the first $n$ terms in such sequences be denoted by $w^{(n)} = (w_0, w_1,$ .......$w_n$), $n \geq 0$. Suppose we know the first $(n-1)$ steps of a sample path — that is, $w^{(n-1)}$ is known. If we assume, on one hand, that the conditional probability of going to state $s_n$ in $t$ or less time units does not depend on the number of transitions $(n-1)$ or the age $x_{n-1}$ at which an individual enters state $s_{n-1}$, then the model is called a homogeneous or age-independent semi-Markov model. If we assume, on the other hand, that the conditional probability of going to state $s_{n-1}$ in $t$ or less time units does not depend on the number of transitions $(n-1)$ but depends on the age $x_{n-1}$ at which an individual enters state $s_{n-1}$, then the model is called an non-homogeneous or age-dependent semi-Markov model. Whether homogeneous or non-homogeneous, the semi-Markov model ignores the number of transitions already made; that is, it ignores how the origin state was reached (the Markovian condition being still valid).

In a semi-Markov model, we have to consider two random variables: $X_n$ denoting the state entered at the $n$-th step and $Y_n$ denoting the sojourn time in state $X_{n-1}$, $n \geq 1$. In order to simplify notations, let $s_n = k$ and $s_{n-1} = j$. Then, we can denote the conditional probability of transition from state $j$ to state $k$ in $t$ or less time units by $A_{jk}(t)$, which is a *direct* one-step transition probability — also called first passage probability. This is a non-negative and non-decreasing function for every pair of states $j$ and $k$ in $S$.

From the first passage probabilities described above, one can derive other results such as duration-stay probabilities (equivalent to survival probabilities), mean-length of stay, and state transition probabilities. Duration-stay probabilities denote the probabilities that an individual entering a specific state at $t=0$ will still remain in the same state after $t>0$ time units. It is the complement of all possible (direct) transitions from a given state and, hence, implies no move whatsoever in the given time interval. The mean-length of stay in a particular state is derived from duration-stay probabilities within a given time-interval by cumulating them over each unit of duration. The state transition probabilities are calculated through first passage probabilities and duration-stay probabilities. These probabilities have been found useful in studies using the semi-Markov model (Mode, 1985; Rajulton, 1988; Rajulton and Balakrishnan, 1990).

A Non-Markov Model

The strategies followed in the Markovian and semi-Markovian frameworks relied on the information on the timing and sequence of transitions. If we are interested in the question "In what way does a sample path traced by one individual differ from the one traced by another in its influence on the event under study?", then we have to consider the third aspect of life history information; namely the number or the order of events. A model that takes into account the history of past transitions becomes non-Markovian.

Including past histories in a life history analysis is worthwhile in spite of the cumbersome procedures involved. No social process obeys the Markovian condition and no social system is without memory. On the contrary, memory of the past pervades individual lives as well as any social system. Even the methodological devices initially built on the Markovian condition have been shown to be influenced by the past in further analysis (for an example in demographic projections, see Gibberd, 1981). There is no doubt, therefore, that the neglect of the past leads to a biased analysis of a social behaviour.

Past history becomes more relevant when we consider the fact that in real situations, transitions among states in one specific system are very often dependent on transitions among states in another system. In other words, reality calls for a simultaneous consideration of two or more dependent subsystems or state spaces rather than one single system or state space. This is because no transitions are made sequentially within one system without being influenced by those in another system. Expressed differently, transitions among one or more systems are dynamically interdependent; transitions in one system are mutually affected by transitions made in the past in another system or in the same system. Demographic examples abound. Transitions from one parity to the next obviously depend on transitions among marital states, and vice versa. This is similar for transitions among states of education and labour force participation, among marital and employment states, and so on.

One of the possible ways of analyzing the mutual dependence of transitions in two or more subsystems is to consider a system of coupled states rather than single states (for details, see Tuma and Hannan, 1984; chapter 9). Let us consider an example of two interdependent subsystems, marital and parity states. We shall denote the two subsystems by $S_M$ and $S_P$, and use a pair of random variables $X_{Mn}$ and $X_{Pn}$ to denote the states occupied at the $n$-th step in each subsystem. These random variables take distinct values denoted by positive integers from 1 to $S_j$, where $S_j$ is the

size of the state space of the *j*-th subsystem, *j=M or P*. Elementary principles of combinatorics suggest that the number of distinct values which the system of coupled states can take is equal to $\Pi_{j=M,P} S_j$. Specifically, if the random variable $X_{Mn}$ denotes marital status taking values as defined earlier, and if the random variable $X_{Pn}$ takes values, say from 0 to 4 to denote the parities, then we have a total of 5*5 = 25 different coupled states $(X_{Mn}, X_{Pn})$ in the system; namely (1,0), (1,1), (1,2), (1,3), (1,4), (2,0), (2,1), (2,2).... and so on. Though there are 25 different coupled states, not all of them would be meaningful, nor would all of them be realizable; only a few would be relevant in practical situations. For example, the first transition from (1,0), that is, Never Married with 0 parity, can be only to either (1,1) or (2,0) or (3,0), representing single parenthood, cohabitation or first marriage, respectively. Similarly, the second transition from, say, (1,1) will be either to (2,1) or (3,1), and so on. What is important is that the whole history is preserved in examining the transitions involved in the newly considered system of coupled states.

In practice, since it is very inconvenient to build models on non-Markovian lines, when attempts are made to include past history, real non-Markovian schemes are always reduced to Markovian (or semi-Markovian) ones. Thus, in order to examine the influence of each subsystem on the other, we can extend the semi-Markovian framework by considering coupled states in any specific order of transition. The algorithm specified for the semi-Markovian model holds good here too, except that now we know how the origin state was reached. Analogous to the definitions used in the case of the semi-Markovian model, we shall consider the triplet $(s_{Mn}, s_{Pn}, y_n)$ as the state occupied by an individual at the *n*-th step after a duration of $y_n$ time units in the previous state. Starting from $w_0=(s_{M0}, s_{P0})$, we follow an individual's history as follows: $w_1=(s_{M1}, s_{P1}, y_1)$, $w_2=(s_{M2}, s_{P2}, y_2)$, and so on. For an illustration of an analysis which considers coupled states in a non-Markovian set-up (but reduced to a semi-Markovian one), see Rajulton and Balakrishnan (1990).

In a non-Markov model, when a specific transition (described by order and sequence) is made only by a handful of individuals, transition probabilities can become erratic and unreliable. The analyst may, therefore, have to impose certain criteria on what is an acceptable error. In order to avoid small numbers, it would be wise to ignore other forms of heterogeneity, such as age at entry into the origin state, region of birth, and other socioeconomic covariates. Limiting the number of states in each subsystem will also be helpful.

Remarks on the Use of Stochastic Processes
The need of stochastic processes in demographic analysis is increasingly being felt, and greater availability of life history information gives hope of satisfying this need. A life history analysis, however, need not be restricted to the above three stochastic frameworks. The above discussions were focussed on the three models not only because of their possible application in many other fields, but also because of their close resemblance to life table techniques, the unifying framework in demographic analysis. A computer program, LIFEHIST, which is usable on a mainframe computer and of these three frameworks to any topic of interest, is available at the Population Studies Centre (PSC) at the University of Western Ontario. The procedures for using the program are explained in a manual which is also available at the PSC.

Obviously, other stochastic processes relevant to life history analysis do exist; for example, diffusion processes which can be very useful in demographic research. Few attempts have been made to explore their usefulness, even in a single state analysis (but see Casterline *et al.,* 1987). Diffusion processes, unlike other models, have distinct advantages which are pointed out in the next section.

Working with stochastic processes turns out a very detailed output, as the whole process is examined from "start" to "finish". Instead of such a detailed analysis, one may resort to a parametric approach to life history analysis (as opposed to "non-parametric" approach that life tables involve). It is possible to parametrize the observed heterogeneity in transition probabilities obtained through the stochastic frameworks; see Mode (1985) and Rajulton (1988) for some examples. However, the direct parametric approach has its own advantages, especially when one wants to consider heterogeneity and further disaggregation of a sample may not be possible. We turn our attention to this point in the next section.

*Hazard Models for Life History Analysis:  Case of Heterogeneity*
In spite of the accumulating literature on the so-called "determinants" of various social behaviours, the earlier mentioned assumption that certain characteristics of individuals and contexts affect the change process has not been explored much. This assumption is closely related to the problem of heterogeneity.

Characteristics (individual or otherwise) that differentiate a population are innumerable. Age is an important factor of heterogeneity and has taken the greatest share in demographic analysis. Apart from age, many other time-invariant individual characteristics such as sex, race and place of birth

as well as time-variant characteristics such as marital, employment and education statuses can be observed. The former have often been used in studying the differential effects on specific behaviour while the latter have not been as well recognized by demographers. Analytical tools that make use of observed characteristics (often expressed as covariates) have been developed and applied in studies that assess the impact of these covariates on human behaviour. The well-known hazard model, in its various forms (Cox, 1972; Vanderhoeft, 1985; Menken *et al.,* 1981; Tuma *et al.,* 1979) is one such analytical tool in survival analysis. Similar hazard models can be used in analyzing life history information, even though it is difficult to see how all the three aspects (number, timing and sequence) can be used without much cost and inconvenience. See Tuma *et al.* (1979), Flinn and Heckman (1982), and Allison (1982), for more details on this specific topic.

While many such characteristics (whether individual or otherwise, whether time-variant or time-invariant) can be observed and measured, many others are neither observable nor measurable. Many unobserved and unobservable characteristics do influence behaviour, and the idea of using stochastic processes to estimate the effects is being explored by a number of analysts.

The awareness that a consideration of unobserved heterogeneity at an individual level would portray reality more adequately led to the suggestion of including the distribution of unobservables in an analysis of observed characteristics. In its simplest forms, this suggestion is not totally new to the field of social science. As early as the 1950s, Silcock (1954) studied job mobility through a model which contained a parametric representation of unobserved heterogeneity in transition rates; he assumed a gamma distribution for the rate of leaving a job. The gamma distribution of hazards lacking measured characteristics has been found to be very useful in many other studies (for examples see Spilerman, 1972; Vaupel *et al.,* 1979). Besides the theoretical reason of flexibility of a gamma distribution to accomodate a variety of situations (especially the unknown and the unobservable ones), it is difficult to see on what other grounds (of explanation, interpretation, and generalization) a gamma distribution should be preferred to other distributions.

Recent research has revealed that an analysis of observed characteristics is to a certain extent subject to the assumption of the kind of distribution of the unobservables (Heckman and Singer, 1982; Yamaguchi, 1986). The study by Heckman and Singer was disturbing to many conscientious researchers; it showed that conclusions can be sensitive to the choice of

distribution of unobserved heterogeneity. Trussell and Richards (1985) followed suit with the finding that conclusions are not only sensitive to the choice of distribution but also to the assumption about the time-dependence in hazard functions.

Another line of research in heterogeneity dynamics (centred at the International Institute of Applied Systems Analysis [IIASA], Laxenburg) showed that a mere assumption of the kind of distribution of unobservables without an explicit consideration of the relationship that could exist between observed rate and parameters of underlying unobserved characteristics would lead to wrong inferences. The random walk model of human mortality proposed by Woodbury and Manton (1977) was developed into a more general model based on conditional gaussian diffusion processes (Yashin *et al.*, 1985). This more general model covers the possibility of analyzing the effects of both the observed and unobserved or partially observed characteristics in their multivariate form. A greater value of the model is that it allows the inclusion of variables that evolve over time; for example, changes in employment, economic conditions, education and marital status are known to have their effects on each other, as was discussed earlier in the section entitled "A Non-Markov Model." The gaussian model's flexibility offers a vast scope for its application to many fields, but needs a practical algorithm as well as empirical verifications.

*Summary*
The review of types of analyses that are possible with life history information covers quite a range of recent innovations in demographic analysis. It can be said unambiguously that in dealing with life history data, an analyst should be willing to address the following complexities of any social process: age dependence, duration dependence, age and duration dependence, systemic relationships, effects of time-varying explanatory variables and inter-dependence of life history events. Not that all these complexities can be addressed simultaneously. One can hardly do so, but certainly it pays to be aware of them. This paper has been intended primarily to bring into focus these various aspects of life history analysis and to suggest viable techniques to address them.

As far as empirical applications are concerned, the above complexities can be progressively addressed as follows:

(1) A Markovian analysis to examine the age-dependence of transitions in a specific system;

(2) A semi-Markovian framework to examine duration-dependence;

(3) An non-homogeneous semi-Markov process to address both age and duration dependence simultaneously;

(4) A parameterization of transition probabilities involved in the semi-Markov process to account for heterogeneity;

(5) A gaussian process to examine the effects of time-varying explanatory variables, both observed and unobserved; and finally

(6) a non-Markovian framework to examine the inter-dependence of life history events.

*References*

Allison, P.D. 1982. Discrete time methods for the analysis of event histories. In S. Leinhardt (ed.), Sociological Methodology 1982. San Francisco, CA: Jossey-Bass.

Bartholomew, D.J. 1967. Stochastic Models for Social Processes. 2nd edn. New York, NY: John Wiley.

Blumen, I., M. Kogan and P.J. McCarthy. 1955. The Industrial Mobility of Labour as Probability Process. Cornell Studies of Industry and Labour Relations, no. 6. Ithaca, NY: Cornell University Press.

Casterline, J.B., M.R. Montgomery and R.L. Clark. 1987. Diffusion models of fertility control: Are there new insights? Working Paper, 87-06. Providence, RI: Population Studies Training Centre, Brown University.

Coleman, J.S. 1964. Introduction to Mathematical Sociology. New York, NY: Free Press.

Cox, D.R. 1972. Regression models and life tables. Journal of the Royal Statistical Society B34:187-220.

Feller, W. 1950. An Introduction to Probability Theory and Its Applications. Vol.1. New York, NY: John Wiley.

_____. 1966. An Introduction to Probability Theory and Its Applications. Vol.2. New York, NY: John Wiley.

Flinn, C. and J.J. Heckman. 1982. New methods for analyzing individual event histories. In S. Leinhardt (ed.), Sociological Methodology 1982. San Francisco, CA: Jossey-Bass

Gibberd, R. 1981. Aggregation of population projection models. In A. Rogers (ed.), Advances in Multiregional Demography. RR-81-6. Laxenburg, Austria: International Institute Applied Systems Analysis.

Ginsberg, R.B. 1971. Semi-Markov processes and mobility. Journal of Mathematical Sociology 1:233-262.

Heckman, J.J. and B. Singer. 1982. Population heterogeneity in demographic models. In K.C. Land and A. Rogers (eds.), Multidimensional Mathematical Demography. New York, NY: Academic Press.

Hoem, J.M. 1972. Inhomogeneous semi-Markov processes, select actuarial tables and duration dependence in demography. In T.N.E. Greville (ed.), Population Dynamics. New York, NY: Academic Press.

Hoem, J.M. and U. Funck Jensen. 1982. Multistate life table methodology: A probabilist critique. In K.C. Land and A. Rogers (eds.), Multidimensional Mathematical Demography. New York, NY: Academic Press.

Howard, R.A. 1971. Dynamic Probabilistic Systems. Vol. 1. New York, NY: John Wiley.

Land, K.C. and A. Rogers (eds.). Multidimensional Mathematical Demography. New York, NY: Academic Press.

Menken, J., J. Trussell, D. Stempel and O. Babakol. 1981. Proportional hazards life table models: An illustrative analysis of socio-demographic influences on marriage dissolution in the United States. Demography 18:181-200.

Mode, C.J. 1982. Increment decrement life tables and semi-Markovian processes from a sample path perspective. In K.C. Land, and A. Rogers (eds.), Multidimensional Mathematical Demography. New York, NY: Academic Press.

_____. 1985. Stochastic Processes in Demography and their Computer Implementation. New York, NY: Springer-Verlag.

Pressat, R. 1969. L'Analyse Demographique. Paris: Presses Universitaires de France.

Rajulton, F. 1988. A semi-Markovian approach to using event history data in multiregional demography. Mathematical Population Studies 1(3):289-315.

Rajulton, F. and T.R. Balakrishnan. 1990. Interdependence of transitions among marital and parity states in Canada. Canadian Studies in Population 17(1):107-132.

Ross, S. 1983. Stochastic Processes. New York, NY: John Wiley.

Schoen, R. and K. Land. 1979. A general algorithm for estimating a Markov-generated increment-decrement life tables with applications to marital status patterns. Journal of the American Statistical Association 74:761-776.

Silcock, H. 1954. The phenomenon of labour turn-over. Journal of the Royal Statistical Society A117:429-440.

Spilerman, S. 1970. The causes of racial disturbances: A comparison of alternative explanations. American Sociological Review 35:627-649.

_____. 1972. The analysis of mobility processes by the introduction of independent variables into a Markov chain. American Sociological Review 37:277-294.

Trussell, J. and T. Richards. 1985. Correcting for unobserved heterogeneity in hazard models using the Heckman-Singer procedure. In N.B. Tuma (ed.), Sociological Methodology. San Francisco, CA: Jossey-Bass.

Tuma, N.B., M.T. Hannan and L.P. Groeneveld. 1979. Dynamic analysis of event histories. American Journal of Sociology 84:820-854.

Tuma, N.B. and M.T. Hannan. 1984. Social Dynamics. New York, NY: Academic Press.

Vanderhoeft, C. 1985. Stratified proportional hazard models: A GLIM-oriented approach with special reference to problem of competing risks. Working Paper no. 5. Brussels, Belgium: Inter-University Program in Demography.

Vaupel, J.W., K.G. Manton and E. Stallard. 1979. The impact of heterogeneity in individual frailty on the dynamics of mortality. Demography 16:434-454.

Willekens, F. and A. Rogers. 1978. Spatial Population Analysis, Methods and Computer Programs. RR-78-18. Laxemburg, Austria: International Institute of Applied Systems Anaysis.

Woodbury, M.A. and K.G. Manton. 1977. A random walk model of human mortality and aging. Theoretical Population Biology 11:37-48.

Yamaguchi, K. 1986. Alternative approaches to unobserved heterogeneity in the analysis of repeatable events. Sociological Methodology 16:213-249.

Yashin, A.I., K.G. Manton and J.W. Vaupel. 1985. Mortality and aging in a heterogeneous population: A stochastic process model with observed and unobserved variables. Theoretical Population Biology 27:154-175.