

Event History Analysis In Demography

Daniel Courgeau
and
Éva Lelièvre

CLARENDON PRESS • OXFORD

1992

Contents

Introduction	1
I. Extending the Scope of Longitudinal Analysis	7
1. Observation of Event Histories	9
1.1. Various types of surveys	11
1.2. The 'Triple Biographie' survey	13
1.3. The problem of weight	16
1.4. Incomplete and erroneous survey observation	18
1.5. The problem of censoring	26
1.6. Conclusion	28
2. Formalization of the Analysis	29
2.1. Analysis of a homogeneous cohort experiencing a single event	30
2.2. Analysis of a heterogeneous cohort and of the interaction between phenomena	36
2.3. Towards a more exhaustive analysis of human behaviour patterns	45
2.4. Conclusion	48
3. Methods of Estimation using Censored Observations	50
3.1. Censoring problems	50
3.2. Right-censoring	52
3.3. Left-censoring	61
3.4. Conclusion	66
4. Study of a Single Event	68
4.1. Single sample: a single event	68
4.2. Single sample: competing risks	72
4.3. Multiple samples: comparative tests	75
4.4. Conclusion	81

5. Reciprocal Study of Interactions between Two Events	82
5.1. Conception of the analysis	82
5.2. The bivariate case	87
5.3. Practical analysis	93
5.4. Conclusion	98
6. Extending to More Complex Situations	99
6.1. Presentation and limits of the practical application	99
6.2. Interactions between three events: two study cases	101
6.3. Interactions between two renewable processes	104
6.4. Conclusion	106
II. Extending the Scope of Regression Models	107
7. Statistical Formalization of Parametric Analysis	109
7.1. Some useful parametric models in demography	109
7.2. Regression models	135
7.3. Conclusion	144
8. Methods of Estimation of Parametric Models	145
8.1. Computation of the likelihood in the presence of censoring	146
8.2. Estimating the parameters and testing their value	148
8.3. Estimation of the parameters: examples	149
8.4. Comparison of parametric models	173
8.5. Conclusion	177
9. Methods of Semi-parametric Analysis	180
9.1. From parametric regressions to semi-parametric proportional hazard models	180
9.2. Methods of estimation	182
9.3. The Newton-Raphson algorithm	187
9.4. The choice of a model for the analysis of interactions	188
9.5. Some applications	189
9.6. Conclusion	193
10. Conclusion	195
10.1. Analysis of interactions between phenomena	195

10.2. Dealing with heterogeneity in populations	200
10.3. New lines of research	204
Appendix	206
Bibliography	210
Index	221

Introduction

Up to the present, longitudinal analysis has basically been developed through an approach that takes each demographic phenomenon into consideration separately. Its main objective has been to isolate each phenomenon in its 'pure state'. In particular, it was considered necessary to separate the effects of each demographic variable from those of the others, such as mortality and migrations. To this end, a certain number of hypotheses were required, but they could not be tested against existing sources (Henry, 1959, 1966).

Longitudinal analysis was developed on the basis of aggregate data such as registration statistics, or, whenever possible, data from population registers. Even though these sources allow each event to be dealt with separately, since they eliminate disturbing effects, they hardly make it possible to analyse interactions between different phenomena. This explains why, in traditional demographic manuals, we find isolated phenomena in their 'pure state' presented as the subject matter of separate chapters: nuptiality, fertility, mortality, moves and migrations (Pressat, 1961; Henry, 1972).

None the less, certain demographers have pointed out the usefulness of analysing interactions between demographic phenomena. Pressat (1966) emphasized that 'the search for correlations between demographic phenomena, even though this domain remains unexplored, should enable us to deepen our knowledge considerably'. He did not, however, give any indication of what method to follow. Similarly, Henry (1972), in his analysis of nuptiality, stated that, 'in the case of out-migrants, one might be tempted to substitute their nuptiality abroad for that which they would have experienced had they stayed in their original country; yet nuptiality in a foreign country depends on conditions that may differ considerably'. This is a definite recognition of the interaction between the two phenomena and the change in an individual's nuptial behaviour following his out-migration. As data were lacking, however, Henry did not pursue this analysis any further.

Another important problem is that of heterogeneity. This point was also touched on by Henry (1959). Although, 'in the case of a homogeneous cohort, the statistical history of the component individuals is the same as the cohort's statistical history', this result no longer holds true when a heterogeneous cohort is being studied. Thus, for example, in the simplest of cases, where two sub-populations each have a constant though different probability of occurrence of the studied event over the period, the overall population will no longer have a constant rate, taken as an average of the two sub-populations' rates. It may well be the case at the outset, but with time the sub-population with the highest hazard rate will be eliminated from the population at risk by a selection process. This means that after a while the hazard rate of the observed population will converge with that of the sub-population with the lowest hazard rate. This heterogeneity may, of course, be of a much more complex nature and would need to be studied further.

A precise knowledge of the practical implications of heterogeneity in human groups would necessitate further differential demographic research into the individual's physical and psychological characteristics, in order to study both the dispersion and correlation of the intra-group demographic indices which have so far been studied in a rather cursory fashion. (Henry, 1959)

As long as demographers use statistics such as those published in registration records or population registers, they have no way of dealing with the two basic problems: the analysis of interactions between demographic phenomena, and the analysis of heterogeneity in human groups. Other sources must now be used in order to observe a group of individuals over their entire lifetime, or at least part of it, as well as to collect a greater number of characteristics for each respondent.

It can thus be seen that the unit of analysis will no longer be the event (death, marriage, birth, migration, etc.): instead, each individual biography will be viewed as a more complex process. The question is no longer one of trying to isolate each phenomenon in its 'pure state': on the contrary, we must try and see how one event in an individual's existence can influence his life-course, and how certain characteristics can induce an individual to adopt behaviour patterns that are different from those of another individual.

This change of view leads us to reformulate the fundamental notions of demographic analysis in terms of complex stochastic processes. Let us take a more detailed look at how this is to be done.

Demographic processes do not occur in an abstract space-time, but originate within a given social structure. Someone born into a Lobi tribe in the early twentieth century will have quite a different biography from someone born into rural France of the same period, or someone born into present-day urban France. In each of these social structures, however, it is possible to distinguish relational systems that have developed to a greater or lesser extent depending on the group or the society concerned: 'family, economic, political, religious, educational, associational and informal systems' (Kimball and Pearsall, 1954). There is, of course, nothing to prevent new types of relational systems from appearing in the future. Our approach does not consider a society as a closed entity, but rather as one in constant evolution.

Each member of a given society is simultaneously involved in the various systems. For example, someone living in present-day France may be part of the family system as a member of an unmarried couple and father of a child; in the economic system as an engineer in the car industry; in the political system as a town councillor; in the religious system as a non-practising Catholic; in the educational system as a recipient of professional training; in the associational system as an amateur footballer; and, finally, in the informal system through his occasional attendance at parent-teacher meetings to solve his child's educational problems.

It is this interaction between the different types of involvement that creates a space and time specific to each situation. The geographical or occupational mobility of a single person may be much more frequent and may occur over greater distances than that of a married person, especially if the latter has one or more children. The married person is naturally tied to his place of residence and work by constraints related to his spouse's work-place, the location of his children's school, etc.

Event history analysis will thus attempt to place these changes in the time and space of an individual's life, in his social context. The point is to see how an event of a family, economic or other nature experienced by the individual will change the probability of other events happening to him over his lifetime. We shall, for instance, try to discover how his marriage can influence his professional career, his spatial mobility and other occurrences, such as the birth of a child or a break with his original family ties.

Here we are directly concerned with the analysis of interactions between demographic phenomena, the utility of which we have

already pointed out. This method of analysis has its place in the study of event histories.

Similarly, when trying to understand an individual's behaviour, one must take into account his social origins and his entire past history. In this case we are supposing that behaviour patterns are not innate but rather that they can change over an individual's lifetime as a result of what he experiences and acquires with time. Thus, two individuals from the same social background, but who have taken entirely different paths in life, can have attitudes to marriage, forming a family, career, etc., that diverge increasingly as time goes on.

We, thus, arrive at a method of analysis of population heterogeneity which uses a dynamic rather than static approach and which, accordingly, has its place in the study of individual event histories.

It is important to note that this analysis of population heterogeneity is not deterministic, but basically probabilistic.

Consequently, a large number of individuals who find themselves in the same conflictual situation at the outset will have different probabilities of finding solutions to the situation before a given date. Some will never find solutions; others may invent for themselves an entirely new pattern of behaviour which can provide a better solution to their conflictual situation. We are, thus, allowing the individual a margin of freedom which may lead to entirely new situations. Such a margin of freedom is of course essential, as no chain of events is predestined, but evolves with the course of time.

On this point, we are very close to Prigogine and Stengers (1988), who recognize that 'the event creates a difference between the past and the future . . . It is the intelligible outcome of a past, from which, however, it could not be deduced. It opens out an historic future where the insignificance or meaning of its consequences will be decided.' Here, one finds that same margin of freedom which can lead to entirely new situations.

After this informal introduction to event history analysis, can we now proceed to describe it in more formal terms?

When an individual is born, his life can follow a wide variety of paths. These different life-courses, however, are far from being equally probable. An individual's event history can therefore be defined as the result of a complex stochastic process, which develops over time, yet is situated within given historical, geographical, economic, and social conditions.

Let Ω_θ represent the set of event histories or partial event histories which can be observed up to a time θ . As already stated, our observation must be limited to the past, since the future brings into play new situations that cannot be deduced from the past. For example, the career of garage-owner could not have been envisaged before the appearance of motor cars, and in the near future genetic discoveries may stretch the span of human life to 150 years or more. The analysis carried out at time t only takes into account past behaviour patterns and projects them into the future without introducing elements affecting their evolution. This evolution, however, takes place at a slow enough rhythm to enable the analysis of the past to enlighten us as to the probability of various events occurring in the near future.

Let ω_n denote an entirely observed event history, where n precedes θ , and ω_θ an event history observed up to the moment θ , which is not finished. It can be said that the events in either of these histories are variables defined in the general space of Ω_θ . For example, the age at which one of these individuals marries in an *application* of Ω_θ on $[0, +\infty]$.

Now, let us give Ω_θ with a sigma-algebra¹ \mathcal{B}_θ of events that are to be analysed within a given population, together with a probability measure P_θ of \mathcal{B}_θ on $[0,1]$, which assigns the different probabilities to the different events occurring within the observed population. In this case, $(\Omega_\theta, \mathcal{B}_\theta, P_\theta)$ does indeed define a probability space.

A random moment T will thus be a function of time on the Ω_θ probability space, which can in this case be extended beyond θ , supposing that the probabilities defined before θ remain the same over time. Thus, for example, if an individual is not married in θ , we may suppose that his age at marriage is a random variable T following the same distribution function as that observed for the individuals already married who have similar characteristics. Besides, this hypothesis is sometimes unnecessary, and one can simply work on event histories that are completely terminated, as is done in historical demography.

The method of event history analysis we present here will thus involve estimating the probability distribution of the life-courses followed by a given population. This distribution may vary from one sub-population to another and may depend on certain characteristics of the individuals in the sub-population (social and economic characteristics of parents or grandparents, for example). These life-courses

¹ The set of the Ω_θ parts.

feature random variables T_1, T_2, \dots, T_n which represent the duration of stay in the different states that constitute them. Of course, these variables are not independent, and the distribution of life-courses to be estimated is the result of their joint distribution.

Our approach, therefore, supposes that individual behaviour can be described as a complex stochastic process. Having once assumed this model, we shall begin with a statistical estimation of the distribution of the variables previously defined which we develop in the present work. Once known, these distributions make it possible to deal with the more complex distribution of the overall life-course.

We shall first examine methods for collecting these event histories. Generally, it is not possible to have an exhaustive collection of biographies. More often, one works with data on partially observed event histories collected from a sample of individuals. It is therefore necessary to examine the different problems raised by such incomplete observation.

We shall then proceed to formalize the methods of analysis and estimation, starting with the most simple case and gradually introducing an increasing degree of complexity. After studying an event, we shall move on to the reciprocal study of the interactions between two events, before extending these non-parametric models to cover more complex situations. Next, we shall examine parametric models which make it possible to introduce the effect of a large number of characteristics on the duration of stay in a given state. Finally, we shall deal with semi-parametric models which combine the two preceding approaches.

These various methods will be illustrated by applications to very different situations, so as to show the possibility for their generalization. In the Appendix, we shall indicate the programmes for carrying out these analyses, so that the reader may actually make use of them.

This book provides both a detailed theoretical presentation of methods for event history analysis and a practical application of these methods to files that exist already or that are to be created on the basis of event history surveys.

Subject Index

- Aalen estimator 59-61, 73
- accelerated failure time model
 - 141-3, 169-73, 180
 - comparison with proportional hazards 142
 - log-logistic 142-3, 169-73
- actuarial estimator:
 - of hazard function 71, 90
- baseline hazard 180-1, 189
- bias 51, 61, 67
- birth 2, 10, 19, 21, 50, 96, 101-2, 109, 196
 - child 2, 29
 - date of 21, 22
 - first 26, 39, 62-5, 85, 102, 193, 199
 - last 19, 52, 87
 - of the last child 85, 113, 151, 153, 155, 159
 - second 11, 85, 86, 199
 - third 86, 193
- bivariate case 82, 87-95, 99, 101-2
- censoring 50-67, 146-7, 181
 - independent 147
 - left 26, 50-1, 61-2
 - right 11, 27, 52-3
 - type of 50
- cohabitation 94, 101
- comparison of distributions 134-5, 173-7
- Competing risks 72-5
- computer packages:
 - EVACOV (INED) 189, 209
 - GLIM 209
 - LIFETEST (SAS) 78-80, 208
 - LOGLIN 208
 - PL2 (BMDP) 189, 209
 - RATE 144, 209
 - ROOT (INED) 95, 209
- confidence interval 48, 92, 150
- consensual union 10, 101
- covariance matrix 54, 76, 91, 156, 162, 164, 169, 188
- cumulative intensity, *see* integrated cumulative hazard
- death 2, 10-11, 20, 29, 50, 54, 56, 82
 - cause of 59, 72
 - date of 21
- density probability function:
 - conditional, *see* hazard function
 - joint 89
 - non-parametric 30-1, 88
 - parametric 110, 114, 118, 122, 126, 128, 129, 131, 137, 138, 142
- dependance 83-7, 90
 - a priori 86-7, 200
 - recipocal 88, 199
 - unilateral 84, 86, 88, 198
- departure 57
 - from the parents' home 27, 83, 85, 93, 94
 - from professional activity 96
 - from the agricultural sector 199
- divorce 94, 101
 - date 9
- duration of stay 6, 16, 19, 23-5, 27, 29, 43, 52, 78-9, 106, 115, 123-4, 166-7, 173, 201, 204
- dwelling 15, 174
 - change of 16-17, 21, 24-5, 115-16, 120, 123-4, 127
 - date of arrival in the 19
 - purchase of the first 191-2
- educational level 135, 155, 157, 162, 173
 - status 109
- emancipation 21-2

- exponential distribution 110–13, 121, 126, 134, 137
 - estimation of the parameter of 149–57, 174
 - mixing 113–21
- family 3, 4, 9, 189, 191
 - history 12–14, 27
- female activity 85, 87, 103–4
- fertility 1, 9, 11, 37–8, 85, 87, 101–2, 104, 199–200
 - hazard rate 97
 - rate 20
 - survey 20
- Fisher information matrix 54, 148, 157, 170
- Fisher–Snedecor distribution 133–4, 175
 - estimation of the parameters of 175–6
- fuzzy time 96–7
- gamma distribution 128–9
- Gompertz distribution 16, 121–5, 127, 134
 - estimation of the parameters of 163–9, 174
- Gompertz–Makeham distribution 124–5, 144
- Greenwood formula 55–6, 69, 71
- hazard function estimation 54, 68, 71, 90
 - non-parametric 31–2, 46, 72, 87, 99, 104–5
 - parametric 110, 114, 118, 121, 125, 128, 129, 131, 137, 138, 143
 - semi-parametric 180
- heterogeneity 2, 4, 30, 36, 109, 117, 135, 195, 197, 200–4
 - unobserved 44–5, 178, 194, 202–3
- home owner 113, 132–2, 139–40, 151, 153, 155, 159, 162–3, 174, 191
- incomplete observation 18–26
- independance 84, 88–9, 198
- instantaneous hazard rate, *see* hazard function
- instantaneous rate of failure, *see* hazard function
- interaction 1, 3, 36, 83, 85–6, 93–4, 95, 101, 104, 137, 189, 203
 - analysis 9–10, 99
- Job:
 - at marriage 85
 - change 29, 52, 102–3
 - first 27, 29, 47, 78–80, 83, 94
 - in farming 84, 189
 - mobility 109, 114
- Kaplan–Meier estimator 53, 55, 68, 73, 186
- likelihood:
 - function 32–7, 47, 54, 145–7, 157
 - log 69, 148, 150
 - partial 32
- log-logistic distribution 131–3, 134–5, 138, 143
 - estimation of the parameters of 169–73, 176
 - with accelerated failure times 142–3
- log-normal distribution 129–31, 176
- Markov 30, 38–45
- marriage 2, 4, 10–11, 14, 21, 24–5, 27, 29, 38, 46, 84–5, 87, 94–5, 101, 109, 189, 191, 193, 199
 - age at 5, 26, 105
 - date of 9, 21–2, 29, 61
 - duration 17, 23
 - first 50
 - status 9
- martingale 60
- maximum likelihood estimator 47, 54, 70–1, 146, 150, 152, 157
- memory 12
 - errors 20–6
- migration 1–2, 26, 29, 38–9, 43, 46, 50, 54, 61, 68, 101, 104, 109, 136, 173–4, 199

- date of 9, 21
- first 26, 65-6, 94, 102
- hazard rate of 77
- history 12-13, 27
- in- 21, 95
- internal 10, 114
- last 19, 52
- multiple 9
- out- 12, 21, 77
- prenuptial 15
- rate 16, 23
- to metropolitan areas 85, 95, 199
- to non-metropolitan areas 95
- mixing distribution:
 - exponential 113-21
- mortality 1, 9, 18, 37, 59, 74-5
 - cumulative hazard curves of 75
 - rate 18
- move 1, 29
 - previous 26
- mover-stayer model 45, 114, 117-120, 123, 127
- Nelson:
 - estimator 60, 73
 - plot 73
- Newton-Raphson algorithm 157-8, 161-2, 165, 170, 176, 187-8
- non-multiplicative hazard model 202
- nuptiality 1, 9, 12, 37-8, 87, 95, 189, 204
 - hazard rate of 95, 190
 - rate 18, 37
- occupation 15
 - change 29, 38, 46
 - last 19
 - professionnal 43
 - status 10, 15, 174
- occupational status 78-81
- order statistic 182
- Pareto distribution 118-19
- Poisson 41-2, 52-3
- population register 1-2, 10
- product integral 34
- product-limit estimator, *see* Kaplan-Meier estimator
- professional career 11, 78
 - life 93, 173
- proportional hazard model 136-43, 165, 180-1, 188, 202
- random loss 146
- rank test 75-6, 91, 182
- residence 43
 - change of 19
 - place of 9, 14
 - region of 39
- sampling plan 17
 - informative 18-19
 - non-informative 14, 17-18
- separation 101-2
- simultaneity 90, 95-7, 100
- spatial mobility 24-5
- survey 11
 - multiround 11
 - prospective 11, 28
 - retrospective 11-13, 26-8
 - 'Triple biographie' 13-17, 20, 48, 61, 63, 77, 113, 115, 151, 162, 166, 173, 189
- survivor function:
 - estimation 55, 68, 71, 185-7
 - non-parametric 2, 30, 88
 - parametric 110, 114, 118, 121, 126, 128, 129, 131, 137, 138, 142
 - semi-parametric 185
- time 4-5, 10, 29-30, 37, 43-8
 - continuous 30, 33-5
 - dependant characteristics 178
 - discrete 32, 33-5
 - interval 41, 50
 - of interview 10
 - of occurrence 46
 - at risk 150
 - of survey 28
 - space- 3
 - waiting 52
- variance 54, 76, 145, 149, 150, 152
 - asymptotic 55-6, 68-9, 71

Weibull distribution 125-7, 134,
138-9
estimation of the parameters of
158-63, 76
weighting 16-18

work:

history 12-16, 27

place of 3

start to 199

Author Index

- Aalen, O. 59, 74, 91
 Allison, P. 36
 Arjas, E. 36

 Blumen, I. 114
 Bretagnolle, J. 203

 Coale, A. 45
 Courgeau, D. 15, 16, 21, 23, 26,
 39, 43, 61, 84-5, 87, 96, 100,
 104, 114, 166, 173, 189, 191,
 199, 201-2, 204
 Cox, D. 90, 148-9
 Crowley, J. 188

 Deroo, M. 20
 Duchêne, J. 10, 21
 Dussaix, A. 20

 Elder, G. 28

 Feller, W. 42, 51
 Firdion, J. 10, 12, 15, 21-2
 Foner, A. 28
 Funck Jensen, U. 40

 Ginsberg, R. 43
 Groot, L. 96-7

 Heckman, J. 178, 203
 Henry, L. 1-2, 37, 195
 Hinkley, D. 148-9
 Hoem, J. 18-19, 40, 91
 Hu, M. 188
 Huber-Carol, C. 203

 Johnson, N. 129

 Kalbfleish, J. 47, 110, 134, 147,
 182
 Kangas, P. 36
 Kaplan, E. 27, 53, 58

 Keilman, N. 96-7
 Keyfitz, N. 38
 Kertzer, D. 28
 Kimball, S. 3
 Klijzing, E. 96-7
 Kotz, S. 129

 Langevin, A. 28
 Lelièvre, E. 16, 84-5, 87, 96, 100,
 104, 189, 191, 199, 201-2, 204
 Lyberg, I. 12, 20

 McCarthy, P. 114
 McGinnis, R. 43
 McNeil, D. 45
 Marvin, K. 114
 Meyer, P. 27, 53, 59
 Modell, J. 28
 Monnier, A. 11
 Murphy, M. 193

 Nelson, W. 59
 Nizard, A. 18

 Oakes, D. 90

 Peto, R. 91
 Pike, M. 91
 Poulain, M. 10, 12, 15, 21-2
 Pourcher, G. 13
 Prentice, R. 47, 110, 134, 147,
 176, 182
 Pressat, R. 1, 195
 Prigogine, I. 4

 Riandey, B. 10, 12, 14-15, 21-2
 Richards, T. 178, 203
 Rogers, A. 36
 Rouy, E. 63

 Schow, G. 92
 Siegers, J. 96-7
 Singer, B. 44, 178, 203

Spillerman, S. 44

Stengers, I. 5

Trussel, J. 178, 203

Turnbull, B. 62

Vaeth, M. 92

Vallin, J. 18

Wendel, B. 10