# Arbeidsnotater

WORKING PAPERS FROM THE CENTRAL BUREAU OF STATISTICS OF NORWAY

IO 70/16                                        Oslo, 21. september 1970

ON THE STATISTICAL THEORY OF

ANALYTIC GRADUATION

By

Jan M. Hoem

C o n t e n t s

1

## PREFACE

The present paper was written while the author visited the
University of California, Berkeley, on a Ford Foundation fellowship. It will
appear in the Proceedings of the Sixth Berkeley Symposium on Mathematical
Statistics and Probability.

# 1. Introduction. Motivating examples.

## 1.1. Example 1: Rates of mortality.

Standard methods for the investigation of human mortality will produce statistics such as those given in extract in table 1. The mortality rate at age $x$ is interpreted as a measure of the mortality risk for women born in the year 1968-$x$, and the corresponding number "exposed to risk" in column 2 is used as a measure of the accuracy of this rate. (This will be made clearer later on.) Unless the population is substantially larger than the one producing these data, the diagram of the sequence of rates, plotted against age, will have a rather rugged appearance. Figure 1, based on the same data as table 1, shows the typical form of such diagrams. There seems to be a universal conviction, however, that "real mortality" would be portrayed by a smooth curve, and that any irregularities of curves of observed mortality rates are due to accidental circumstances. The observed rates are then regarded as "raw" or primary estimates of the underlying "real" rates, and graduation is employed to get a smoother curve.

A number of techniques have been developed to graduate age-specific mortality rates, as can be seen from any text on the subject. (See, for instance, [54], [58, pp. 145-197], [81, pp. 216-237, 243-244, and 251-252].) Most of these methods have been developed by intuitive arguments, at least initially, but investigations of statistical properties of some of them have also appeared [1], [2], [43], [44], [46], [60], [68], [70], [74], [81, p. 252]. One class of such methods consists in fitting a parametric function to the observed rates. We shall call this the class of analytic graduation methods.

Quite a number of functions have been suggested for analytic graduation of mortality rates [45, pp. 236-238], [66, pp. 453-454], [77, pp. 79-85], [81, pp. 56-60 and 243-244]. By far the most commonly used for the adult ages is the Gompertz-Makeham formula

$$(1.1) \quad g_x(\alpha, \beta, c) = \alpha + \beta c^x \quad \text{for } \beta > 0, \ c > 1, \ \alpha > - \beta c^{x_{min}},$$

where $x$ represents age attained. We have fitted this function to our data in figure 1 by minimum $\chi^2$. Other common methods are least squares and some moments methods. We shall describe each of these in turn.

## Table 1

### Age - specific mortality

### Females, municipality of Oslo, Norway, 1968

| Age (1) | Exposed to risk[*] (2) | Deaths[**] (3) | Mortality rate per thousand[+] (4) |
|---|---|---|---|
| 40 | 2798 | 1 | 0.357 |
| 41 | 2924.5 | 3 | 1.025 |
| 42 | 3156 | 6 | 1.901 |
| 43 | 3272.5 | 6 | 1.833 |
| 44 | 3465.5 | 6 | 1.731 |
| 45 | 3639 | 11 | 3.022 |
| 46 | 3770 | 5 | 1.326 |
| 47 | 4057 | 10 | 2.464 |
| 48 | 3886.5 | 10 | 2.573 |
| 49 | 3650.5 | 10 | 2.739 |
| .. | ...... | .. | ..... |
| .. | ...... | .. | ..... |
| 80 | 1204.5 | 111 | 92.154 |
| 81 | 1064 | 81 | 76.127 |
| 82 | 930 | 118 | 126.881 |
| 83 | 835 | 90 | 107.784 |
| 84 | 709.5 | 101 | 142.353 |
| 85 | 615.5 | 74 | 120.227 |
| 86 | 502 | 97 | 193.227 |
| 87 | 408 | 68 | 166.666 |
| 88 | 341 | 66 | 193.548 |
| 89 | 378 | 53 | 140.211 |
| 90 | 218.5 | 43 | 196.796 |

Source:  Central Bureau of Statistics of Norway.

[*] Arithmetic mean of the number of persons at a given age as of January 1, 1963, and the corresponding number as of December 31, 1963.

[**] Age at death is taken as 1968 minus year of birth.

[+] Ratio between entries in columns (3) and (2), multiplied by 1000.

1.2 Example 2:  Rates of fertility.  A standard
investigation of age-specific human fertility will pro-
duce a table quite similar to table 1, except of course
that column (3) there will contain numbers of births
(or usually numbers of liveborn children) by age of
mother.  A corresponding diagram will look something
like the one in figure 2, and graduation will again
give a smoother curve.

A fertility curve of this sort closely resembles
certain density functions, and one category of functions
proposed for the analytic graduation of fertility curves
consists of densities from the Pearson family [13], [27],
[45, pp. 140-169], [52], [55], [71], [75], [78], [79],
in particular Pearson type I, III, IV, VI, and the nor-
mal density, multiplied by a constant.

Another category of graduating functions consists
of polynomials in x [11], [27], like

(1.2)  $b_x(a,b,c,d) = (x - \alpha + 1)(\beta - x)^2 (a + bx + cx^2 + dx^3)$,

where x stands for age of mother at childbearing, and
where $[\alpha,\beta)$ is the fertile period for females.  It is
customary to take $\alpha = 15$ and $\beta = 45$ or $\beta = 50$, but in
certain cases $\alpha$ and $\beta$ occur as parameters which are
estimated [13], [53], [75].

The Hadwiger function,

(1.3)  $h_x(R,T,H,d) = \dfrac{RH}{T\sqrt{\pi}} (\dfrac{T}{x-d})^{3/2} \exp\{-H^2(\dfrac{T}{x-d} + \dfrac{x-d}{T} - 2)\}$,

with $R > 0$, $T > 0$, $H > 0$, $d < \alpha$, is a third type of graduating formula [27], [28], [31], [45, pp. 149-169], [78], [79]. (We follow Yntema's notation.) Other functions have also been suggested [12], [47], [48], [49], [53], [70a].

Naturally, the same type of functions will be used for the graduation of other vital rates whose diagrams have the same general form as fertility rates, such as marriage rates [22, pp. 99-101], [48].

1.3 Introductory. The above age-specific rates
of mortality and fertility are examples of the kind of
vital rates which occur in fields like actuarial science,
biostatistics, and demography. In the present paper we
shall give a contribution to the statistical theory of
curve-fitting as applied to such rates in general. We
shall suggest a probabilistic model within which the
rates appear as estimators for certain parameters called
forces of transition, and shall show how the analytic
graduation can be interpreted as a procedure used to
further estimate a set of "more basic" parameters, viz.
those of the graduating function.

The model will be introduced in section 2. In
sections 3 and 4, we describe how the rates appear
within the model, and sections 5 to 9 are devoted to the
study of analytic graduation methods. We shall be con-
cerned mainly with the asymptotic statistical properties
(as the population size N increases) of the estimators.
Most of our results are straightforward consequences of
general asymptotic theory, and we shall often use
standard theorems from that field, like those of chapters
4 and 5 in [10] and Theorem 4.2.5 in [3], without ex-
plicit reference. We shall quote references whenever
we use a deeper result.

Since we use standard theorems, it is not sur-
prising that we can prove theorems which correspond to

previous standard results. Thus (speaking informally here) we shall see that none of the general estimation procedures we study will be better than one of the maximum likelihood type, and that a (modified) minimum $\chi^2$ procedure is equally good, while moment methods will usually give less favourable results.

We feel that there may be a need for some explanation why procedures of the type which we shall describe are preferred to certain others. Rather than breaking up our presentation of the techniques involved by giving parts of this explanation as we go along, we have preferred to include it all in section 10.

Apart from what is contained in sections 1.1 and 1.2 above, no numerical examples will be given in this paper. Numerical investigations are planned and will be reported at a later date.

## 2. A Markov process model.

### 2.1. The general model.
To describe the phenomena in hand we shall use a Markov process model. Let $y_t$ be the sample function value at time $t$ of a time-inhomogeneous Markov process with a denumerable state space I and a continuous time parameter restricted to some finite time interval $[0,\zeta)$. Let the transition probabilities be

$$P_{iJ}(s,t) = P\{y_t \in J \mid y_s = i\},$$

for $0 \leq s < t < \zeta$, $i \in I$, $J \subseteq I$, and assume that

$P_{ii}(s,t) \equiv 1$, $\lim P_{ij}(s,t) \equiv \delta_{ij}$ (a Kronecker delta)

as $t \downarrow s$. We introduce the <u>forces of transition</u>,

$$\mu_{ij}(s) = \lim_{t \downarrow s} P_{ij}(s,t)/(t-s) \quad \text{for } i \neq j,$$

and the <u>forces of decrement</u>,

$$\mu_i(s) = \lim_{t \downarrow s} \{1 - P_{ii}(s,t)\}/(t-s),$$

for $0 \leq s < \zeta$, and assume that all $\mu_i$ and $\mu_{ij}$ are finite and integrable over $[0,\zeta)$. We also assume that

$$(2.1) \qquad \mu_i = \sum_{j \in I-i} \mu_{ij} \quad \text{for each } i \in I.$$

We shall call a state $i$ absorbing if $\mu_i = 0$.

The problem which leads us to study analytic graduation consists in finding a method for estimating one or more of the $\mu_{ij}(\cdot)$ from data of the type which one encounters within the fields of application mentioned at the beginning of subsection 1.3.

2.2 Examples. We shall give some examples to show how models in the applications appear as particular cases of the general model in section 2.1 above.

(1) Our simplest example will be a model with only two states, called "alive" (state 1) and "dead" (state 2). State 2 is absorbing, and there is only one non-zero force of transition and of decrement, viz.

$$\mu(\cdot) = \mu_1(\cdot) = \mu_{12}(\cdot),$$

called <u>the force of mortality</u>. The rates of section <u>1.1</u> will be seen to appear within this model. The time parameter is represented by a person's age.

(<u>ii</u>) The age-specific fertility rates of section <u>1.2</u> can be interpreted within a model with a double infinity of states. A woman will be said to be in state (k,1) at age x if (she is alive then and) her parity is k, i.e., she has had k births, k = 0, 1, 2, ⋯ . She will be said to be in state (k,2) at age x if she has died within age x and her parity at death was k. All states of the form (k,2) are absorbing. We select two suitable functions, $\mu(\cdot)$ and $\varphi(\cdot)$, and set

$$\mu_{(k,1),(k,2)}(\cdot) = \mu(\cdot),$$

and

$$\mu_{(k,1),(k+1,1)}(\cdot) = \varphi(\cdot),$$

for k = 0, 1, ⋯ , while all other $\mu_{ij} = 0$. The function $\mu$ will be called the force of mortality, and $\varphi$ will be called <u>the force of fertility</u>. Again the time parameter is represented by the woman's age. This model, which we have studied in some detail previously [36], is not particularly "realistic", but it is probably the simplest one in which the rates of section <u>1.2</u> can be meaningfully discussed. More realistic fertility models of this type have appeared elsewhere [37], [38].

(iii) To describe <u>marriage formation and dissolution</u>, we suggest a model with five states, called "never married" (state 1), "married" (state 2), "widowed" (state 3), "divorced" (state 4), and "dead" (state 5). State 5 is absorbing. The following forces correspond to impossible direct transitions, and are therefore identically equal to zero: $\mu_{11}$ for $i > 1$, $\mu_{34}$, $\mu_{43}$, $\mu_{5j}$ for $j < 5$. The model applies to one sex only, while the other sex appears only implicitly, as a kind of shadow factor. We have also looked at marriage models elsewhere [41].

Other models of this type have been studied by, for instance, DuPasquier [21], Sverdrup [69], Simonsen [65], Chiang [15, chapters 4, 5, and 7], and Hoem [35]. Compare also [25] and [63].

In each of these models, a state $i \in I$ corresponds to some vital status, i.e., a marital status, a social status, a birth parity, and so on. A transition similarly corresponds to a vital event, such as a death, a birth, a marriage, a divorce, and so on. An individual sample path will be visualized intuitively as a person (or sometimes a group of persons, like a household or a family) moving through some of the statuses of the system specified. The sample paths will be taken as stochastically independent.

2.3 Seniority. In demographic models one often wishes to distinguish between an age parameter (which may be actual age obtained, duration of marriage, interval since last previous birth, etc.), calendar time, and observational time. It is the age parameter which corresponds to the time parameter in the Markov process of section 2.1. In a general model it may be useful to have a separate name for this (unspecified) age parameter, covering all interpretations which it may have in the applications. Following Henry [33] who calls it ancienneté, we shall use the name seniority for it.

2.4 Some basic assumptions about the forces of transition. In what follows, we shall disregard the forces of transition which are identically equal to zero because they correspond to impossible direct transitions by the definition of the model. Even if the state space I can be (countably) infinite, there are many cases where only a finite number of the non-zero forces of transition are distinct. (Compare example 2.2 (ii).) We shall assume everywhere that there exist A non-negative real functions, $\lambda_1, \cdots, \lambda_A$, such that each $\mu_{ij}$, not identically zero, equals some $\lambda_a$, and such that for each $\lambda_a$ there exists a $\mu_{ij} = \lambda_a$. By (2.1), we may then write

(2.2) $\qquad \mu_i = \sum\limits_{a=1}^{A} c_a(i) \lambda_a$ for each $i \in I$,

where

(2.3) $\qquad c_a(i) = \sum\limits_{\{j \in I-i : \mu_{ij}=\lambda_a\}} 1.$

Each $c_a(i) < \infty$ because all $\mu_i < \infty$. (2.2) shows that for any given $i \in I$, exactly $\sum_a c_a(i)$ of the $\mu_{ij}$ can be positive, i.e., a finite number of the $\mu_{ij}$ only, while the rest are identically equal to zero.

Let us also assume everywhere that

(2.4) $\qquad \sup\{\mu_i(s) : 0 \le s < \zeta, i \in I\} < \infty.$

This assumption is not necessary for what follows, and it can be relaxed [39], [40, section 5], but one will expect it to hold in practice and it will simplify our exposition. It follows from (2.4) [40, section 3.1] that there only exists a finite number of <u>distinct</u> vectors $\underset{\sim}{c}(i) = (c_1(i), \cdots, c_A(i))$. Let us call these $(c_{b1}, \cdots, c_{bA})$ for $b = 1, 2, \cdots, B$, and let

(2.5) $\qquad \gamma_b = \sum\limits_{a=1}^{A} c_{ba} \lambda_a$ for $b = 1, 2, \cdots, B.$

Then each nonzero $\mu_i$ equals some $\gamma_b$ and for each $\gamma_b$ there is a $\mu_i$ which equals it. There is, thus, only a finite number of distinct forces of decrement as well. (This need not be the case if (2.4) does not hold.) For each $i \in I$, let $b(i)$ be defined by $\mu_i = \gamma_{b(i)}$. Then, by (2.3),

$$(2.6) \qquad c_{b(1),a} = \sum_{\{j \in I-1 : \mu_{1j}=\lambda_a\}} 1 .$$

Let

$$\underset{\sim}{c} = \begin{pmatrix} c_{11}, & \cdots & , & c_{1A} \\ \cdot & \cdot \cdot \cdot \cdot \cdot & \cdot \\ c_{B1}, & \cdots & , & c_{BA} \end{pmatrix} .$$

We shall finally assume everywhere that the rank of $\underset{\sim}{c}$ is B. The extension to rank $\underset{\sim}{c} < B$ is easy [40, section 3.1], but it leads to slightly more complicated formulas.

### 3. The primary or "raw" estimates.

### 3.1. Approximation of the $\lambda_a$ by step functions.

As a first step in our description of the kind of estimation methods which have produced the rates of section 1.1 and 1.2, we shall approximate the $\lambda_a$ by step functions. The seniority interval $[0,\zeta)$ is partitioned into D subintervals, $[\zeta_0,\zeta_1)$, $[\zeta_1,\zeta_2)$, $\cdots$ , $[\zeta_{D-1},\zeta_D)$, with $\zeta_0 = 0$ and $\zeta_D = \zeta$. Let $I_d(\cdot)$ be the indicator function of the interval $[\zeta_{d-1},\zeta_d)$, and let

$$(3.1) \qquad \lambda_a^{\#}(\cdot) = \sum_{d=1}^{D} \lambda_{ad} I_d(\cdot).$$

Here each $\lambda_{ad}$ is a constant chosen in such a way that it can represent the values of $\lambda_a$ in $[\zeta_{d-1}, \zeta_d)$.

If $\lambda_a$ is assumed to be a nice and smooth function, with certain known monotonicity properties, say, then $\lambda_a^{\#}$ will inherit these properties, modified, of course, by the fact that the latter is a step function.

In what follows, we shall assume that the $\lambda_a^{\#}$ give an adequate representation of the $\lambda_a$, and our calculations will be made as if we actually had $\lambda_a = \lambda_a^{\#}$ for $a = 1, 2, \cdots, A$.

3.2 More about the $\zeta_d$. In this presentation, we use the same partitioning $\{\zeta_d: \ d = 0, 1, \cdots, D\}$ for all $\lambda_a$. In certain situations one would rather use different partitionings for different $\lambda_a$. The results of this paper will continue to hold for such cases with only quite obvious modifications [37].

The approach sketched in section 3.1 is closely related to histogram methods for the estimation of a probability density or a generalized failure rate [6], [72], [73]. Whereas the lengths of the histogram intervals are often made to converge to zero as the number of observations increase, however, this is not the case for the seniority subintervals above. The $\zeta_d$ will typically be selected according to conventional rules established with different considerations in mind than statistical convergence properties. When $\zeta$ is of the size order of several decades, as is often the case,

the seniority interval $[0,\zeta)$ will usually be partitioned into one-year or five-year intervals, possibly with a longer "tail" interval at the upper end. There is a tendency to use shorter subintervals in a large population than in a small one, at least if the data are reliable, but an interval length shorter than one year is commonly used only in certain standardized contexts, such as in investigations of infant mortality, where $\zeta$ equals one year of age [8, p. 211], [30, tables 38 to 42], [66, p. 84]. There seems to be no affinity to the idea of letting such interval lengths decrease to zero.

3.3 On the observational plan. There are a number of observational plans (or ascertainment methods) in use in the fields of application we have in mind, and one could construct others from ideas used in life testing. (See, e.g., [20], [32], [64], [80].)

In the present paper we shall only consider observational plans where a group of people are followed continuously over some time interval $[0,T)$. The data collection will consist of noting what happens to each person while under observation, i.e., which states of $I$ he visits and just when vital events occur to him.

It is characteristic for the types of populations which occur in practice that some people enter them and others leave during the study period. They may also be heterogeneous with respect to seniority, in the sense that

those who come under observation (whether they are in
the population from the outset or enter later on) may
have different seniorities at time  T.  We want to
cover such possibilities.  Let, therefore,  N  be the
number of individuals ever getting observed.  Let us
say that person no. k  enters the population at some
time  $t_k \epsilon [0,T\rangle$  with seniority  $x_k$  and a status cor-
responding to state  $r_k$,  and that he stays there at
least until time  $t_k + z_k \epsilon [0,T]$, when observation is
discontinued.  We shall take the entrance time  $t_k$,
the initial seniority  $x_k$,  the initial state  $r_k$,  and
the exposure time  $z_k$  to be preassigned, i.e., not
random.  (Other possibilities are discussed in [40].)
Any period spent in an absorbing state, for instance
after the death of the individual, is included in the
period of exposure  $[t_k, t_k + z_k\rangle$,  although of course
no actual observation is made after a path has entered
such a state.

   We shall also take  N  to be non-random.

   **3.4.  Estimation of the  $\lambda_{ad}$.**  We get an es-
timator for  $\lambda_a^\#$  by plugging estimators  $\hat{\lambda}_{ad}$  for the
$\lambda_{ad}$  into the right hand side of (3.1).  (This is how
the rugged curves in figures 1 and 2 have arisen.)
Standard estimators used for this purpose are occur-
ence/exposure rates, like in sections 1.1 and 1.2
[61], [62].  We shall see how these arise.

Some of the $\lambda_{ad}$ may be known to be zero because they correspond to vital events which are impossible during $[\zeta_{d-1}, \zeta_d)$, such as births after menopause. We will take the other $\lambda_{ad}$ to be strictly positive. Let

$$\mathcal{G} = \{(a,d): \lambda_{ad} > 0\}.$$

We can regard $\{\lambda_{ad}: (a,d)\epsilon\mathcal{G}\}$ as a point in the space

$$\Lambda_0 = \underset{(a,d)\epsilon\mathcal{G}}{\times} \{x_{ad} > 0\}.$$

The situation in hand will usually restrict the possible points we actually can have to a proper subset $\Lambda$ of $\Lambda_0$. We shall take $\Lambda$ to be open.

Now let $M_k(a,d)$ be the number of transitions observed for path no. k during the seniority interval $[\zeta_{d-1}, \zeta_d)$, direct from any state i to any state j where $\mu_{ij} = \lambda_a$. Let $U_k(i,d)$ be the total time spent in state i during $[\zeta_{d-1}, \zeta_d)$ by this path, and let

$$(3.2) \quad V_k(b,d) = \underset{\{i\epsilon I: \mu_i = \gamma_b\}}{\Sigma} U_k(i,d).$$

Then $V_k(b,d)$ is the total time spent in any state i where $\mu_i = \gamma_b$ by path k during the interval mentioned. Finally, let

$$(3.3) \quad \begin{aligned} L_k(a,d) &= \overset{B}{\underset{b=1}{\Sigma}} c_{ba} V_k(b,d) = \\ &= \underset{i\epsilon I}{\Sigma} c_{b(i),a} U_k(i,d), \end{aligned}$$

and let us use the notation

$$X = \sum_k X_k,$$

where $X_k$ is any quantity depending on $k$. Since the forces of transition are represented by step functions, we can then use the same method as in [37, §4.8] to write the likelihood in the form

$$\prod_{(a,d)\in G} \lambda_{ad}^{M(a,d)} \cdot \exp\{-\sum_{d=1}^{D}\sum_{b=1}^{B}\gamma_{bd}V(b,d)\}$$

$$= \prod_{(a,d)\in G} \lambda_{ad}^{M(a,d)} \cdot \exp\left\{-\sum_{(a,d)\in G}\lambda_{ad}L(a,d)\right\},$$

where

$$\gamma_{bd} = \sum_{a=1}^{A}c_{ba}\lambda_{ad}.$$

(Compare (2.5).) Thus we are dealing with a Darmois-Koopman class of probability distributions, and one may show [40, section 3.1] that $\{M(a,d), V(b,d): a = 1, 2, \cdots, A; b = 1, 2, \cdots, B; d = 1, 2, \cdots, D\}$ is minimal sufficient for the $\lambda_{ad}$. An unrestricted maximization of the likelihood function would give the estimators

(3.4)     $\hat{\lambda}_{ad} = M(a,d)/L(a,d),$

which are the occurrence/exposure rates we mentioned. (We arbitrarily set $\hat{\lambda}_{ad} = 0$ if $L(a,d) = 0$.) The point $\{\hat{\lambda}_{ad}:(a,d)\in G\}$ need not lie in $\Lambda$, nor even in $\Lambda_0$ if some $\hat{\lambda}_{ad} = 0$. However, under certain conditions,

spelt out in theorem 1 below, the probability that the point lies in $\Lambda$ increases to 1 as $N \to \infty$

3.5. Asymptotic properties of the $\hat{\lambda}_{ad}$. For each $(a,d) \in \mathfrak{C}$, the variables $L_1(a,d), \cdots, L_N(a,d)$ will not generally be identically distributed unless all $(x_k, z_k, r_k)$ are equal. Similarly for $M_1(a,d), \cdots, M_N(a,d)$. Nevertheless one can prove the following consistency theorem [40, section 4.2]:

Theorem 1: Assume that $P\{L(a,d) > 0\} \to 1$ as $N \to \infty$, and that a finite positive limit

$$(3.5) \qquad L_{ad} = \lim_{N \to \infty} EL(a,d)/N$$

exists. Then $\hat{\lambda}_{ad}$ converges to $\lambda_{ad}$ in probability as $N \to \infty$.

To arrive at a theorem concerning the asymptotic distribution of $\hat{\lambda}_{ad}$ as $N \to \infty$, we make an additional set of assumptions, which establish a grouping of the $(x_k, z_k, r_k)$ at a finite set of strategic values. More precisely we make

Assumption 1: There exists a finite set of possible initial seniorities, $y_1, \cdots, y_H$, a finite set of possible exposure times, $w_1, \cdots, w_J$, and a finite set of possible initial states, $s_1, \cdots, s_Q$, such that each $(x_k, z_k, r_k)$ must equal some $(y_h, w_j, s_q)$. We let

$$K_{hjq} = \{k: (x_k, z_k, r_k) = (y_h, w_j, s_q)\},$$

and let $S_{hjq}(N)$ be the number of elements in $K_{hjq}$. We assume that

$$\alpha_{hjq} = \lim_{N \to \infty} S_{hjq}(N)/N$$

exists for each $(h, j, q)$.

If $\varepsilon_{hjq}(a,d) = EL_k(a,d)$ for $k \in K_{hjq}$, we get under assumption 1 that the $L_{ad}$ of (3.5) satisfy

(3.6)     $$L_{ad} = \sum_{h,j,q} \alpha_{hjq} \, \varepsilon_{hjq}(a,d).$$

We may then prove ([40, section 4.2]; compare [46a])

Theorem 2: Under assumption 1, the variables $N^{\frac{1}{2}}(\hat{\lambda}_{ad} - \lambda_{ad})$ for which $P\{L(a,d) > 0\} \to 1$ as $N \to \infty$, are asymptotically independent and normally distributed with means $\underline{0}$ and asymptotic variances

(3.7)     $$\sigma_{ad}^2 = \text{as} \cdot \text{var } N^{\frac{1}{2}}(\hat{\lambda}_{ad} - \lambda_{ad}) = \lambda_{ad}/L_{ad}.$$

We note that under the assumptions of theorem 2,

(3.8)     $$\hat{\sigma}_{ad}^2 = N \, \hat{\lambda}_{ad}/L(a,d)$$

is a consistent estimator for $\sigma_{ad}^2$. Thus we see a justification of the use of the number exposed to risk

(which is  L(a,d)  here) as an intuitive measure of
the accuracy of the corresponding rate of transition
$\hat{\lambda}_{ad}$, as mentioned in section 1.1 for a special case.

We also note that we do not need to know the
value of  N  in order to estimate the  $\lambda_{ad}$  and the
asymptotic variance  $\sigma^2_{ad}/N$  of the  $\hat{\lambda}_{ad}$.

## 4.  Non-observation of part of the state space.

4.1.  A problem.  In section 3 we assumed that
one could observe what state a sample path visited
at any time.  This need not be the case in practice.
Let us give two examples.

(i)  Demographic studies will often be con-
cerned with people living in a restricted area, like
a country or part of a country, and there will be some
in- and out-migration.  Say that a study of marriages
is carried out, perhaps based on a model like the one
in section 2.2(iii).  If a person intially lives in
the study area, then leaves and stays away for a while,
and subsequently returns while the study is still
being conducted, it rarely happens that his changes
of marital status (if any) while outside the study
area are traced.  In many cases one will know his
marital status on departure from the study area, as
well as his status as he returns, but nothing more.

(ii) Similar problems occur in studies of the mortality of insured lives. A person may cancel his insurance policy and be uninsured for a while, then take out a new policy, which may be cancelled again after a while, and so on. The insurer will keep track of deaths among the persons covered by his policies, but will not usually know what happens to the uninsured.

The question is how one should take account of phenomena like these in the estimation procedures.

### 4.2. Formalization of the two examples.

(i) To describe example (i) above in terms of a probabilistic model, let $I_1 = \{1,2,3,4,5\}$ be the state space of example 2.2(iii), let $I_2 = \{1,2\}$, and let $I = I_1 \times I_2$. An individual with marital status $j$ will be said to be in state $(j,1)$ if he lives in the study area, and in state $(j,2)$ if he lives outside it. A migration out of the study area will correspond to a transition from a state $(j_1,1)$ to a state $(j_2,2)$. A migration into the study area will correspond to a transition from a state $(j_2,2)$ to a state $(j_1,1)$. In most cases, $j_1 = j_2$. In any case, we shall take $j_1$ to be observable. Whatever moves the sample path otherwise makes while in the subspace $\{(j,2): j \epsilon I_1\}$ will not be observed.

(ii) To formalize the second example above,

let us use four states, called "alive and insured"
(state 1), "alive and uninsured" (state 2), "dead
while insured" (state 3), and "dead while uninsured"
(state 4). Which transitions are possible and which
are not follows directly from the state names. Except
for transitions from state 2 to state 4, all transi-
tions (and the dates on which they occur) are recorded.
(This is essentially the model studied by DuPasquier
[21], Fix and Neyman [25], and Sverdrup [69], except
that they took all transitions as recorded. Recording
problems different from the present one have been studied
by Høyland [42], Kruopis [46a], and others.)

Let us take all forces of transition to be con-
stants. This will suffice for our purposes, which are
those of illustration. Generalization to other cases
is simple. We shall take the forces of mortality of
the insured and the uninsured to be equal, and let
$\mu = \mu_{13} = \mu_{24}$. Let $\nu = \mu_{12}$, $\rho = \mu_{21}$, $\alpha = \mu + \nu$, $\beta = \mu + \rho$.
Sample path no. k is followed over the period $[0, z_k]$, and
we say that $k \in \kappa$ if this sample path is in state 2 or 4
at time $z_k$, i.e., if person no. k is uninsured then .
All N paths start in state 1, and they make a total
number of $M_{ij}$ jumps from state i to state j, for
$(i,j) \in \{(1,2),(1,3), (2,1)\}$. Let W denote the total
time spent in state 2 by the paths $k \notin \kappa$, and let V
be the total time spent in state 1 by all paths taken

together. For $k \in K$, let $z_k - U_k$ be the time of the last jump recorded from state 1 to state 2 for path $k$, i.e., the time to last observed cancellation. Then the corresponding likelihood can be written as

$$e^{-\alpha V - \beta W} \nu^{M_{12}} \rho^{M_{21}} \mu^{M_{13}} \beta^{-K} \prod_{k \in K} (\mu + \rho e^{-\beta U_k}),$$

where $K$ is the number of elements in $K$. The maximum likelihood estimator of $\nu$ turns out to be

$$\hat{\nu} = M_{12}/V,$$

which is what (3.4) would have given. ($M_{12}$ is the number of cancellations observed.) Closed, explicit expressions for the maximum likelihood estimators of $\mu$ and $\rho$ do not exist in this case. We can still get an estimator of $\mu$, however, by letting

$$\hat{\mu} = M_{13}/V.$$

($M_{13}$ is the number of insured deaths.) The properties of $\hat{\nu}$ and $\hat{\mu}$ will appear by specialization of the results in section 4.3 below.

4.3. The general case. Consider now the general model with the assumptions made in sections 2 and 3.1. Let the state space I be partitioned into two disjoint subsets, H and J, and assume that all transitions between states in H can be recorded, while no transitions between states in J are recorded. Any transition

from a state in H to one in J is recorded, as is
also all jumps from J to H. For both kinds of jumps,
one also records the state <u>to</u> which the jump is made.

We redefine the quantities $M(a,d)$ and $L(a,d)$,
initially introduced in section <u>3.4</u>, as follows:

Let G be the set of the a for which there
exists a $\mu_{ij}$, with $i \in H$, such that $\mu_{ij} = \lambda_a$.
For each $a \in G$ and each $d \in \{1, 2, \cdots, D\}$ let $M(a,d)$
now be the total number of transitions observed during
the seniority interval $[\varsigma_{d-1}, \varsigma_d)$, for all paths taken
together, direct from any state $i \in H$ to any state
$j \in I-i$ such that $\mu_{ij} = \lambda_a$. Furthermore, let

$$(4.1) \qquad L(a,d) = \sum_{i \in H} c_{b(i),a} \, U(i,d),$$

with $c_{b(i),a}$ given by (2.6) and $U(i,d)$ defined as
in section <u>3.4</u>. For $a \in G$, let $\hat{\lambda}_{ad}$ be given by
(3.4) with the <u>new</u> definitions of $M(a,d)$ and $L(a,d)$.
<u>Then theorems 1 and 2 hold verbatim for the $a \in G$</u>,
even though the $\hat{\lambda}_{ad}$ need not be maximum likelihood
estimators, as demonstrated in example (<u>ii</u>) above.
If there does not exist any $\mu_{ij}$, with $i \in J$, such
that $\mu_{ij} = \lambda_a$ for any $a \in G$, the $\hat{\lambda}_{ad}$ <u>will</u> be
maximum likelihood estimators, in the sense that they
maximize the likelihood under free variation of the
$\lambda_{ad}$ in $\Lambda_0$.

If the state $j$ cannot be recorded when there is a jump from a state $i \in H$ to a state $j \in J$, the results above continue to hold, provided we again redefine the quantities involved in a natural way. In the definition of $M(a,d)$, we must only include jumps from $i$ to $j$ where both $i$ and $j \neq i$ belong to $H$, and where $\mu_{ij} = \lambda_a$. $G$ is similarly reduced. This time we also redefine $c_{ba}$ by letting

$$(4.2) \quad c_{b(i),a} = \sum_{\{j \in H-i : \mu_{ij} = \lambda_a\}} 1 \quad \text{for } i \in H, \ a \in G.$$

Using (4.2), we define $L(a,d)$ for $a \in G$ by (4.1).

## 5. Conventions and notation relating to analytic graduation

### 5.1. Analytic graduation.

Although an original $\lambda_a$ is assumed to be a nice and smooth function, the estimators $\hat{\lambda}_{ad}$ now in use, such as those in (3.4), will typically produce a $\hat{\lambda}_a^{\#}$ which is considered too irregular, except in large populations. (Compare the account on page 561 in [18].) Analytic graduation then consists in selecting some nice, parametric function $g_a(\cdot, \underset{\sim}{\theta}_a)$ and some representative senority $\xi_d$ from each interval $[\zeta_{d-1}, \zeta_d)$, and in getting an estimator $\hat{\underset{\sim}{\theta}}_a$ for $\underset{\sim}{\theta}_a$ by fitting the values $\{g_a(\xi_d, \underset{\sim}{\theta}_a): d = 1, 2, \cdots, D\}$ to $\{\hat{\lambda}_{ad}: d = 1, 2, \cdots, D\}$ by a suitable method. The function $g_a(\cdot, \hat{\underset{\sim}{\theta}}_a)$, usually regarded as a function of a <u>continuous</u> seniority variable x, represents the final estimator for the function $\lambda_a(\cdot)$.

Most methods for constructing an estimator $\hat{\underset{\sim}{\theta}}_a$ are based on analogies with estimation methods used in other contexts [1], [58], [81]. We shall study least squares and minimum $\chi^2$ methods in section $\underset{\sim}{6}$ [7], [12], [16], [27], [28], [47], [48], [49], [60]. (See also [59].) In section $\underset{\sim}{7}$, we shall discuss moment methods [11], [13], [28], [45, pp. 140-169], [52], [55], [71], [75], [78], [79], and in

section 8 we shall introduce a technique of the maximum likelihood type. Some authors have also used methods involving the minimization of sums of absolute deviations [17], [28].

5.2. Further assumptions and conventions. We shall be working with a single, fixed value of a, and shall therefore suppress this subscript except where it may cause confusion.

In what follows, we shall disregard the fact that some of the $\lambda_d$ may be known to equal zero. The case where some $\lambda_d$ actually do equal zero needs only trivial notational modifications.

Let

$$g_d(\underset{\sim}{\theta}) = g(\xi_d, \underset{\sim}{\theta}),$$

$$\underset{\sim}{g}(\underset{\sim}{\theta}) = (g_1(\underset{\sim}{\theta}), \cdots, g_D(\underset{\sim}{\theta}))'.$$

(The prime denotes a transpose.)

Assumption 2: We assume that $\underset{\sim}{\theta}$ varies in an open subset $\Theta$ of the G-dimensional Euclidean space $\underset{\sim}{R}_G$, where $G < D$. Let $\underset{\sim}{g}$ be a one-to-one, bicontinuous, continuously differentiable mapping of $\Theta$ into

$$\Lambda_0 = \underset{d=1}{\overset{D}{\times}} \{x_d > 0\}.$$

**Define**

$$
\underset{\sim}{J}(\underset{\sim}{\theta}) = \begin{pmatrix} \frac{\partial}{\partial\theta_1} g_1(\underset{\sim}{\theta}) \,, \;\cdots, \; \frac{\partial}{\partial\theta_G} g_1(\underset{\sim}{\theta}) \\ \cdot \;\cdot\;\cdot\;\cdot\;\cdot\;\cdot\;\cdot\;\cdot\;\cdot\;\cdot\;\cdot\;\cdot \\ \frac{\partial}{\partial\theta_1} g_D(\underset{\sim}{\theta} \,, \cdots, \; \frac{\partial}{\partial\theta_G} g_D(\underset{\sim}{\theta}) \end{pmatrix},
$$

<u>and assume that</u> $\underset{\sim}{J}(\underset{\sim}{\theta})$ <u>has rank</u> $G$ <u>for each</u> $\underset{\sim}{\theta} \in \Theta$.

We denote the true value of $\underset{\sim}{\theta}$ by $\underset{\sim}{\theta}^0$, and let

$\underset{\sim}{J}_0 = \underset{\sim}{J}(\underset{\sim}{\theta}^0)$, $\underset{\sim}{\lambda}^0 = \underset{\sim}{g}(\underset{\sim}{\theta}^0)$, and $L_d^0 = \lim\limits_{N\to\infty} E_{\underset{\sim}{\theta}^0} L(d)/N$.

(Compare (3.5.).) We also let

$$
\sigma_{d0}^2 = \lambda_d^0/L_d^0 \quad \text{and} \quad \underset{\sim}{\Sigma}^0 = \mathrm{diag}(\sigma_{10}^2, \;\cdots, \; \sigma_{D0}^2)
$$

(compare (3.7)), with the convention that we write

$\underset{\sim}{M} = \mathrm{diag}(m_1, \;\cdots, \; m_S)$ if $\underset{\sim}{M}$ is a diagonal $S \times S$-

matrix with the $m_s$ as diagonal elements.

Let us denote it by a right topscript $N$ if we

want to stress that a quantity depends on $N$.

In sections 3 and 4 we brought out some estimators

$\hat{\underset{\sim}{\lambda}}^{(N)} = (\hat{\lambda}_1^{(N)}, \;\cdots, \; \hat{\lambda}_D^{(N)})'$ of the common

occurrence/exposure type for the parameter $\underset{\sim}{\lambda} =$

$(\lambda_1, \;\cdots, \; \lambda_D)'$, and we stated some theorems concerning

their asymptotic properties. In much of what follows,

it is precisely these properties which are of interest,

and not the form of the estimators themselves. In

sections 6 and 7, therefore, we shall take $\hat{\underset{\sim}{\lambda}}^{(N)}$ to be

any estimator for $\underset{\sim}{\lambda}$, not necessarily the one given

by (3.4), and we shall continuously make

Assumption 3.

$$N^{\frac{1}{2}}(\hat{\underset{\sim}{\lambda}}(N) - \underset{\sim}{\lambda}^0) \xrightarrow[\underset{\sim}{\theta}^0]{\mathcal{L}^0} \mathfrak{N}(\underset{\sim}{0}, \underset{\sim}{\Sigma}_0),$$

where $\mathfrak{N}(\underset{\sim}{0}, \underset{\sim}{\Sigma}_0)$ is the multinormal distribution with
mean $\underset{\sim}{0}$ and a positive definite covariance matrix
$\underset{\sim}{\Sigma}_0$, which need not be the same as $\underset{\sim}{\Sigma}^0$.

## 6. Analytic graduation through minimization of

## a quadratic form.

6.1. The graduation method. Let $\underset{\sim}{M}$ be a

positive definite, symmetric $D \times D$ matrix whose ele-

ments $m_{ij}$ may (but need not) be random variables. Let

$$Q(\underset{\sim}{\theta}) = N(\hat{\underset{\sim}{\lambda}} - \underset{\sim}{g}(\underset{\sim}{\theta}))' \underset{\sim}{M}(\hat{\underset{\sim}{\lambda}} - \underset{\sim}{g}(\underset{\sim}{\theta})).$$

Assume that there exists a $\underset{\sim}{\theta}$, say $\hat{\underset{\sim}{\theta}}$, which minimizes

$Q(\underset{\sim}{\theta})$. We shall then take $\hat{\underset{\sim}{\theta}}$ to be our estimator for $\underset{\sim}{\theta}$.

A whole class of graduation methods is generated

by the various choices of the matrix $\underset{\sim}{M}$. Thus if we

take $\underset{\sim}{M} = \underset{\sim}{I}$, the identity matrix, we get

$$Q(\underset{\sim}{\theta}) = N \sum_{d=1}^{D} (\hat{\lambda}_d - g_d(\underset{\sim}{\theta}))^2,$$

and $\hat{\underset{\sim}{\theta}}$ becomes a least squares estimator. An analogy

with the modified minimum $\chi^2$ method results from setting

(6.1) $\qquad \underset{\sim}{M} = \text{diag}(1/\hat{\sigma}_1^2, \cdots, 1/\hat{\sigma}_D^2),$

where the $\hat{\sigma}_d^2$ are given by (3.8). We then get

$$Q(\underset{\sim}{\theta}) = \sum_{d=1}^{D} \{M(d)-L(d)g_d(\underset{\sim}{\theta})\}^2/M(d).$$

If, in particular, $\underset{\sim}{g}(\underset{\sim}{\theta})$ is a linear function of $\underset{\sim}{\theta}$, say

(6.2) $\qquad \underset{\sim}{g}(\underset{\sim}{\theta}) = \underset{\sim}{J}_0 \, \underset{\sim}{\theta} + \underset{\sim}{g}_0$

where $\underset{\sim}{J}_0$ is a known $D \times G$ matrix of rank $G$ and $\underset{\sim}{g}_0$ is a known $D \times 1$ vector, we get

$$\underset{\sim}{\hat{\theta}} = (\underset{\sim}{J}_0{}'\, \underset{\sim}{M}\, \underset{\sim}{J}_0)^{-1} \, \underset{\sim}{J}_0{}' \, \underset{\sim}{M}(\underset{\sim}{\hat{\lambda}}-\underset{\sim}{g}_0).$$

A particular case of (6.2) is given in (1.2).

6.2. Asymptotic theory. Let $\{\underset{\sim}{M}^{(N)}\}$ be a sequence of positive definite, symmetric, possibly random, $D \times D$-matrices. For simplicity we assume that the $\underset{\sim}{M}^{(N)}$ are not functions of $\underset{\sim}{\theta}$. (This can be modified. Compare, e.g., [14, Theorem 5].) Let $\underset{\sim}{\hat{\theta}}^{(N)}$ be a value of $\underset{\sim}{\theta}$, if any, which minimizes $Q(\underset{\sim}{\theta})$ with $\underset{\sim}{M} = \underset{\sim}{M}^{(N)}$ and $\underset{\sim}{\hat{\lambda}} = \underset{\sim}{\hat{\lambda}}^{(N)}$. We can then prove the following theorem by the methods of general asymptotic statistical theory. (See [50]. All the hard parts of the proof can be handled by the argument in [9].)

Theorem 3: Make assumptions 2 and 3, and assume also that

$$\text{plim } \underset{\sim}{M}^{(N)} = \underset{\sim}{M}_0,$$

where $\underset{\sim}{M}_0$ is positive definite. With a probability increasing to 1 as $N \to \infty$, there then exists a value $\hat{\underset{\sim}{\theta}}^{(N)} \in \Theta$ which minimizes $Q(\underset{\sim}{\theta})$, and

$$N^{\frac{1}{2}}(\hat{\underset{\sim}{\theta}}^{(N)} - \underset{\sim}{\theta}^0) \xrightarrow{\mathcal{L}_{\theta^0}} \mathfrak{N}(\underset{\sim}{0}, \underset{\sim}{\Sigma}),$$

where

$$(6.3) \quad \underset{\sim}{\Sigma} = (\underset{\sim}{J}_0' \underset{\sim}{M}_0 \underset{\sim}{J}_0)^{-1} \underset{\sim}{J}_0' \underset{\sim}{M}_0 \underset{\sim}{\Sigma}_0 \underset{\sim}{M}_0 \underset{\sim}{J}_0 (\underset{\sim}{J}_0' \underset{\sim}{M}_0 \underset{\sim}{J}_0)^{-1}$$

is positive definite.

Corollary: $$N^{\frac{1}{2}}\{g(\hat{\underset{\sim}{\theta}}^{(N)}) - \underset{\sim}{\lambda}^0\} \xrightarrow{\mathcal{L}_{\theta^0}} \mathfrak{N}(\underset{\sim}{0}, \underset{\sim}{J}'_0 \underset{\sim}{\Sigma} \underset{\sim}{J}_0).$$

Remark 1: If $\underset{\sim}{M}_0 = \underset{\sim}{\Sigma}_0^{-1}$, as is the case when we use (3.4) and (6.1), we get $\underset{\sim}{\Sigma}$ equal to

$$(6.4) \quad \underset{\sim}{\Sigma}_{00} = (\underset{\sim}{J}'_0 \underset{\sim}{\Sigma}_0^{-1} \underset{\sim}{J}_0)^{-1}.$$

Remark 2: Since $G < D$, $\underset{\sim}{J}'_0 \underset{\sim}{\Sigma} \underset{\sim}{J}_0$ is singular.

Remark 3: If we regard $\hat{\underset{\sim}{\theta}}^{(N)}$ as a mapping from $\underset{\sim}{R}_D$ to $\underset{\sim}{R}_G$ (i.e., a function of $\hat{\underset{\sim}{\lambda}}$), we obviously have

$$\hat{\underset{\sim}{\theta}}^{(N)}(g(\underset{\sim}{\theta})) = \underset{\sim}{\theta} \quad \text{for } \underset{\sim}{\theta} \in \Theta,$$

for any positive definite $\underset{\sim}{M}^{(N)}$.

6.3 The choice of $\{M^{(N)}\}$. Since different sequences $\{M^{(N)}\}$ give rise to estimators $\{\hat{\theta}^{(N)}\}$ which may have different asymptotic covariance matrices, one will want to know how to select a $\{M^{(N)}\}$ so as to get a $\Sigma$ which is as favourable as possible. Given two such matrices, $\Sigma_1$ and $\Sigma_2$, where $\Sigma_2 - \Sigma_1$ is positive semidefinite, we shall regard $\Sigma_1$ as the more favourable, since each of the variances on its diagonal will be no greater than the corresponding variance on the diagonal of $\Sigma_2$. At the same time, $J_0' \Sigma_1 J_0$ will be preferred to $J_0' \Sigma_2 J_0$ (compare the corollary to theorem 3), since also $J_0'(\Sigma_2 - \Sigma_1)J_0$ will be positive semidefinite. The following theorem tells us that an $\{M^{(N)}\}$ with $M_0 = \Sigma_0^{-1}$ will be optimal in this sense.

Theorem 4: Let $\Sigma$ and $\Sigma_{00}$ be given by (6.3) and (6.4), respectively. Then $\Sigma - \Sigma_{00}$ is positive semidefinite under the assumptions of theorem 3.

Proof: (i) Let $A$ be any $D \times G$ matrix of rank $G < D$. Then $A(A'A)^{-1}A'$ is idempotent, so all its characteristic roots equal 0 or 1. Thus $I - A(A'A)^{-1}A'$ has only 0 and 1 as characteristic roots, and this matrix, therefore, is positive semidefinite.

(ii) Let us then prove that $\Sigma_0 - J_0(J_0' \Sigma_0^{-1} J_0)^{-1} J_0'$ is positive definite. Let $B$ be a nonsingular matrix such that $B' \Sigma_0 B = I$. Let $v$ be an arbitrary $D \times 1$-vector, and let $w = B^{-1} v$. Then

$$v'\{\Sigma_0 - J_0(J_0'\Sigma_0^{-1} J_0)^{-1} J_0'\}v = w'\{I - B'J_0(J_0'BB'J_0)^{-1} J_0' B\}w$$

$$= w'\{I - A(A'A)^{-1} A'\} w,$$

with $A = B' J_0$. Our assertion then follows from step (i) above.

(iii) Finally, let $v$ be as above, and let $w = M_0 J_0(J_0' M_0 J_0)^{-1} v$. Then $J_0' w = v$ and so

$$v'\{\Sigma - \Sigma_{00}\}v = w'\{\Sigma_0 - J_0 \Sigma_{00} J_0'\} w \geq 0$$

by step (ii) above. Thus $\Sigma - \Sigma_{00}$ is positive semidefinite. $\square$

6.4 The choice of $\{\hat{\underset{\sim}{\lambda}}{}^{(N)}\}$. In sections 6.1 to 6.3 above, we have focused on a single estimator $\{\hat{\underset{\sim}{\lambda}}{}^{(N)}\}$. Assume now that two such sequences are proposed, say $\{\hat{\underset{\sim}{\lambda}}_1^{(N)}\}$ and $\{\hat{\underset{\sim}{\lambda}}_2^{(N)}\}$, both satisfying the assumptions of theorem 3, with asymptotic covariance matrices $\frac{1}{N}\underset{\sim}{\Sigma}_1$ and $\frac{1}{N}\underset{\sim}{\Sigma}_2$, respectively. Say that $\underset{\sim}{\Sigma}_2 - \underset{\sim}{\Sigma}_1$ is positive semidefinite. Intuitively one would expect $\{\hat{\underset{\sim}{\theta}}_1^{(N)}\}$ to have a more favourable asymptotic covariance matrix than $\{\hat{\underset{\sim}{\theta}}_2^{(N)}\}$, when $\{\hat{\underset{\sim}{\theta}}_i^{(N)}\}$, (for i = 1,2) is produced from $\{\hat{\underset{\sim}{\lambda}}_i^{(N)}\}$ by the method of section 6.1 with a choice of $\{\underset{\sim}{M}_i^{(N)}$ which is optimal according to theorem 4. This turns out to be correct:

Theorem 5: If $\underset{\sim}{\Sigma}_1$ and $\underset{\sim}{\Sigma}_2$ are positive definite and $\underset{\sim}{\Sigma}_2 - \underset{\sim}{\Sigma}_1$ is positive semidefinite, then $\underset{\sim}{\Sigma}_{02} - \underset{\sim}{\Sigma}_{01}$ is positive semidefinite, where

$$\underset{\sim}{\Sigma}_{0i} = (\underset{\sim}{J}_0' \; \underset{\sim}{\Sigma}_i^{-1} \; \underset{\sim}{J}_0)^{-1} \qquad \text{for i = 1, 2.}$$

Here $\underset{\sim}{J}_0$ is any $D \times G$ matrix of rank $G < D$.

Proof: If $\underset{\sim}{A}$ and $\underset{\sim}{B}$ are positive definite $D \times D$-matrices with positive semidefinite $\underset{\sim}{A} - \underset{\sim}{B}$, then $\underset{\sim}{B}^{-1} - \underset{\sim}{A}^{-1}$ will also be positive semidefinite [26, page 55, theorem 2.5]. From this the theorem easily follows. []

## 7. Moment methods.

### 7.1. The graduation method.

A moment method estimator $\tilde{\theta}^{(N)}$ of $\theta$ is defined as a solution of the system of equations

$$(7.1) \qquad \sum_{d=1}^{D} \xi_d^r \{\hat{\lambda}_d^{(N)} - g_d(\tilde{\theta}^{(N)})\} = 0 \quad \text{for } r = 0,1, \cdots, G-1,$$

if it exists. Let

$$(7.2) \qquad \underset{\sim}{M} = \begin{pmatrix} 1, & 1, & \cdots, & 1 \\ \xi_1, & \xi_2, & \cdots, & \xi_D \\ \xi_1^2, & \xi_2^2, & \cdots, & \xi_D^2 \\ \xi_1^{G-1}, & \xi_2^{G-1}, & \cdots, & \xi_D^{G-1} \end{pmatrix} .$$

Then (7.1) can be rewritten as

$$(7.3) \qquad \underset{\sim}{M} \{\hat{\underset{\sim}{\lambda}}^{(N)} - \underset{\sim}{g}(\tilde{\underset{\sim}{\theta}}^{(N)})\} = 0 .$$

We shall extend this definition, and shall call $\tilde{\underset{\sim}{\theta}}^{(N)}$ a generalized moment method estimator for $\theta$ if it is a solution of (7.3), where $\underset{\sim}{M}$ here can be any $G \times D$ matrix, i.e., $\underset{\sim}{M}$ need not be given by (7.2).

To give an example of an estimator generated by (7.3) but not satisfying (7.1), we shall consider the

King-Hardy method of estimating the three parameters, $\alpha$, $\beta$, and $c$, of the Gompertz-Makeham function in (1.1). Say that we can take $[0, \varsigma)$ to be the age interval $[x_0, x_0+3h)$ for some integer $h$, and that $\varsigma_x = x_0+x$ for $x = 0, 1, \cdots, 3h$, so that we have one-year age intervals. Then the King-Hardy estimators are the solution $(\tilde{\alpha}, \tilde{\beta}, \tilde{c})$ of the equations

$$(7.4) \qquad \sum_{x=x_0+(k-1)h}^{x_0+kh-1} (\tilde{\alpha}+\tilde{\beta}\,\tilde{c}^x) = \tilde{H}_k \quad \text{for} \quad k = 1, 2, 3,$$

where

$$\tilde{H}_k = \sum_{x=x_0+(k-1)h}^{x_0+kh-1} \hat{\lambda}_x.$$

We get [58, p. 167]

$$(7.5) \qquad \tilde{c}^h = (\tilde{H}_3-\tilde{H}_2)/(\tilde{H}_2-\tilde{H}_1),$$

and similar formulas for $\tilde{\alpha}$ and $\tilde{\beta}$. If we let $\underset{\sim}{m}_x$ be a $1 \times h$ vector where all elements equal $x$, and let

$$\underset{\sim}{M} = \begin{pmatrix} \underset{\sim}{m}_1 & \underset{\sim}{m}_0 & \underset{\sim}{m}_0 \\ \underset{\sim}{m}_0 & \underset{\sim}{m}_1 & \underset{\sim}{m}_0 \\ \underset{\sim}{m}_0 & \underset{\sim}{m}_0 & \underset{\sim}{m}_1 \end{pmatrix},$$

then (7.3) reduces to (7.4) in the case where $g_x(\alpha, \beta, c) = \alpha + \beta\, c^{(x_0+x)}$.

In applications to analytic graduation, the matrix $\underset{\sim}{M}$ is usually non-random and not a function of $N$ or $\underset{\sim}{\theta}$. For simplicity we shall only study this case, but generalization to possibly random $\underset{\sim}{M}$, possibly depending on $N$ and $\underset{\sim}{\theta}$, can be made by standard methods [24], [76].

If, in particular, $\underset{\sim}{g}(\underset{\sim}{\theta})$ is given by (6.2), we get

$$(7.6) \qquad \underset{\sim}{\tilde{\theta}}^{(N)} = (\underset{\sim}{M}\, \underset{\sim}{J}_0)^{-1}\, \underset{\sim}{M}(\underset{\sim}{\hat{\lambda}}^{(N)} - \underset{\sim}{g}_0),$$

provided $\underset{\sim}{M}\, \underset{\sim}{J}_0$ is nonsingular.

### 7.2. Asymptotic theory.

**Theorem 6:** <u>Make assumptions 2 and 3, and assume also that $\underset{\sim}{M}\, \underset{\sim}{J}(\theta)$ is nonsingular for any $\underset{\sim}{\theta} \in \Theta$. There then exists a neighbourhood $\Omega$ of $\underset{\sim}{g}(\Theta)$ and a one-to-one mapping $\underset{\sim}{\tilde{\theta}}^{(N)}$ from $\underset{\sim}{R}_D$ to $\underset{\sim}{R}_G$, continuous in $\Omega$,</u> such that

$$\underset{\sim}{\tilde{\theta}}^{(N)}(\underset{\sim}{g}(\underset{\sim}{\theta})) = \underset{\sim}{\theta} \qquad \text{for } \underset{\sim}{\theta} \in \Theta$$

and (7.3) holds for all $\underset{\sim}{\hat{\lambda}}^{(N)} \in \Omega$. Furthermore,

$$N^{\frac{1}{2}}(\underset{\sim}{\tilde{\theta}}^{(N)} - \underset{\sim}{\theta}^0) \xrightarrow{\mathcal{L}_{\theta}^0} \mathfrak{N}(0, \underset{\sim}{\Sigma}),$$

where

$$(7.7) \qquad \underset{\sim}{\Sigma} = (\underset{\sim}{M}\, \underset{\sim}{J}_0)^{-1}\, \underset{\sim}{M}\, \underset{\sim}{\Sigma}_0 \{(\underset{\sim}{M}\, \underset{\sim}{J}_0)^{-1}\, \underset{\sim}{M}\}'.$$

$\underset{\sim}{\Sigma} - \underset{\sim}{\Sigma}_{00}$ if positive semi-definite. ($\underset{\sim}{\Sigma}_{00}$ is given in (6.4).

Proof: By theorem 1 in [24, p. 1054] we need
only prove the final assertion above. Let $\underset{\sim}{v}$ be an
arbitrary $G \times 1$-vector, and let $\underset{\sim}{w} = (\underset{\sim 0}{J'} \underset{\sim}{M'})^{-1} \underset{\sim}{v}$. Then

$$\underset{\sim}{v}' \{\Sigma - \underset{\sim 00}{\Sigma}\} \underset{\sim}{v} = (\underset{\sim}{M} \underset{\sim}{w})' \{\underset{\sim 0}{\Sigma} - \underset{\sim 0}{J} \underset{\sim 00}{\Sigma} \underset{\sim 0}{J}'\} (\underset{\sim}{M} \underset{\sim}{w}) \geq 0$$

by step (ii) of the proof of theorem 4. $\square$

Remark 1: By the final assertion of the theorem, the
generalized moment method can never give a more favour-
able asymptotic covariance matrix for the estimator of
$\underset{\sim}{\theta}$ than the corresponding "optimal" estimator found in
section 6.

Remark 2: Since $\overset{\wedge}{\underset{\sim}{\lambda}}(N)$ is $N^{\frac{1}{2}}$-consistent for $\underset{\sim}{\lambda}^0$,

$$P\{\overset{\wedge}{\underset{\sim}{\lambda}}(N) \in \Omega\} \to 1 \text{ as } N \to \infty.$$

Remark 3: The analogues of theorem 5 and the corollary
to theorem 3 hold in the present situation.


When the generalized moment method is applied
to a particular case, it is frequently modified to suit
the characteristics of the situation in hand. We shall
give examples of this in sections 7.3 and 7.4.

7.3. Modifications, example 1: Gompertz-Makeham gradua-
tion. In mortality studies using the Gompertz-Makeham

formula (1.1), one will frequently find that $c$ is estimated by (7.5), but that estimators for $\alpha$ and $\beta$ are subsequently found by some other method, for instance by minimizing

$$\sum_{x=x_0}^{x_0+3h-1} (\hat{\lambda}_x^{(N)} - \alpha - \beta \, \tilde{c}^x)^2$$

[81, p. 225]. Let us consider a slightly more general case, and let us estimate $\alpha$ and $\beta$ by minimizing

$$Q(\alpha, \beta) = N(\hat{\underset{\sim}{\lambda}}^{(N)} - \alpha \underset{\sim}{e} - \beta \, \tilde{\underset{\sim}{\psi}}^{(N)})' \, \underset{\sim}{M}^{(N)} (\hat{\underset{\sim}{\lambda}}^{(N)} - \alpha \underset{\sim}{e} - \beta \tilde{\underset{\sim}{\psi}}^{(N)}).$$

Here $\{\underset{\sim}{M}^{(N)}\}$ is a sequence of matrices of the kind studies in section 6, $\underset{\sim}{e}$ is a $3h \times 1$-vector where all elements equal 1, and

$$\tilde{\underset{\sim}{\psi}}^{(N)} = (\tilde{c}^{x_0}, \tilde{c}^{x_0+1}, \cdots, \tilde{c}^{x_0+3h-1})'.$$

Assuming that $\underset{\sim}{M}^{(N)}$ is positive definite, and letting

$$\underset{\sim}{K}^{(N)} = (\underset{\sim}{e}, \tilde{\underset{\sim}{\psi}}^{(N)}),$$

we get the estimators

$$\begin{matrix} \hat{\hat{\alpha}} \\ \hat{\hat{\beta}} \end{matrix} = (\underset{\sim}{K}^{(N)'} \, \underset{\sim}{M}^{(N)} \, \underset{\sim}{K}^{(N)})^{-1} \, \underset{\sim}{K}^{(N)'} \, \underset{\sim}{M}^{(N)} \, \hat{\underset{\sim}{\lambda}}^{(N)}.$$

Now let $\underset{\sim}{m}_x$ be defined as below (7.5), let $\alpha_0, \beta_0$, and $c_0$ be the true values of the parameters, let

$$\underset{\sim}{\ell}_0 = (c_0^{x_0}, \cdots, c_0^{x_0+3h-1})', \quad \underset{\sim}{K}_0 = (\underset{\sim}{e}, \underset{\sim}{\ell}_0),$$

$$\gamma = (c_0-1)/\{h\beta_0 c_0^{x_0+h-1}(c_0^h - 1)^2\},$$

and

$$(7.8) \qquad \underset{\sim}{\Phi} = \begin{pmatrix} (\underset{\sim}{K_0'} \ \underset{\sim}{M}_0 \ \underset{\sim}{K}_0)^{-1} \ \underset{\sim}{K}_0 \ \underset{\sim}{M}_0 \\ \\ (\underset{\sim}{m}_0, \ \underset{\sim}{m}_{-\gamma}, \ \underset{\sim}{m}_\gamma) \end{pmatrix} .$$

We then get

Theorem 7: Let plim $\underset{\sim}{M}^{(N)} = \underset{\sim}{M}_0$, where $\underset{\sim}{M}_0$ is is positive definite, and make assumption 3. Then

$$N^{\frac{1}{2}}\{(\hat{\alpha},\hat{\beta},\tilde{c})' - (\alpha_0,\beta_0,c_0)'\} \xrightarrow{\mathcal{L}(\alpha_0, \beta_0, c_0)} \mathfrak{n}(\underset{\sim}{0}, \underset{\sim}{\Phi} \ \underset{\sim}{\Sigma}_0 \underset{\sim}{\Phi}') .$$

Stevens [67], Patterson [56], Lipton and McGilchrist [51],and some of their references have studied the estimation of the Gompertz-Makeham parameters in a regression model. Stevens [67] found that King-Hardy's method may be very inefficient there. The estimators developed for the regression model can also be used for purposes of analytic graduation, and it would be interesting to see an investigation of their merits in that context.

### 7.5  Modifications, example 2:  Hadwiger

graduation.  Consider now the problem of graduating a

set $\{\hat{\lambda}_x : x = \alpha, \alpha + 1, \cdots, \beta - 1\}$ of female fertility

rates calculated for single-year age groups by fitting

the Hadwiger function (1.3) to the rates.  If we regard

$h_x$ as a function of a continuous $x$, and define

$$R_k'(R,T,H,d) = \int_d^\infty x^k \, h_x(R,T,H,d) \, dx,$$

then

(7.9)  $R_0'(R,T,H,d) = R, \quad R_1'(R,T,H,d) = R \cdot (T + d),$

and the formulas for $R_k'$ for $k \geq 2$ can be found from

the fact that the corresponding cumulants for $d = 0$ are

$$\mathcal{H}_k = 1 \cdot 3 \cdot \cdots (2k-3) 2^{k-1} H^{2k-2} / T^{3k-4} \quad \text{for } k \geq 2$$

(Compare [45, pp. 150, 151, 160].)  No such nice formulas

are known for the discrete case, i.e., for

$$R_k(R,T,H,d) = \sum_{x=\alpha}^{\beta-1} x^k h_x(R,T,H,d),$$

where only integer values of $x$ are used in the summation.

Rather than attempting cumbersome calculations with the

$R_k$, and acting on the analogy between the $R_k$ and the

$R_k'$, Yntema [28], [79] has suggested an estimation pro-

cedure which amounts to the following:

Regard $h_x$ as a function of a <u>continuous</u> x.

Let $U = T + d$.

Then

(7.10) $\qquad h_U(R,T,H,d) = \dfrac{RH}{T\sqrt{\pi}}$ ,

and the mode of the function is

(7.11) $\qquad M = d + 3T\{(1 + 16H^4/9)^{\frac{1}{2}} - 1\}/(4H^2)$.

One easily sees that $M < U$. Solve (7.11) with respect to H, introduce the result into (7.10), let

$$a = \frac{4}{3} \pi(Th_U/R)^2, \quad b = (M-d)/T,$$

and get

$$b = \{(1 + a^2)^{\frac{1}{2}} - 1\}/a.$$

For the range of values in which a will usually lie, the right hand side here is approximately equal to $1 - a^{-1}$. Solving $b \approx 1 - a^{-1}$

with respect to T after substituting $U - T$ for d, we get

(7.12) $\qquad T \approx R^2/\{\frac{4}{3} \pi (U-M)h_U^2\}$ .

Now introduce the estimators as follows: Let

$$\hat{R} = \sum_{x=\alpha}^{\beta-1} \hat{\lambda}_x, \quad \hat{U} = \sum_{x=\alpha}^{\beta-1} x\,\hat{\lambda}_x/\hat{R} .$$

(Compare (7.9).) If [y] denotes the integer value of y, let

$$V = [\hat{U}+\tfrac{1}{2}], \quad \hat{h} = \hat{\lambda}_V, \quad \hat{M} = \min\{x: \hat{\lambda}_x \geq \hat{\lambda}_y \text{ for all } y\}.$$

Finally, let

$$\hat{T} = \hat{R}^2/\{\tfrac{4}{3} \pi(\hat{U} - \hat{M})\hat{h}^2\}, \quad \hat{d} = \hat{U} - \hat{T},$$

(compare (7.12)), and let

$$\hat{H} = \hat{h}\,\hat{T}\,\sqrt{\pi}/\hat{R} = \tfrac{3}{4}\,\hat{R}/\{\sqrt{\pi}\,\hat{h}(\hat{U} - \hat{M})\} \ .$$

(Compare (7.10).) Then $\hat{\underset{\sim}{\theta}} = (\hat{R},\hat{T},\hat{H},\hat{d})$ is an estimator for $\underset{\sim}{\theta} = (R,T,H,d)$. To study its asymptotic properties, we let $\underset{\sim}{\theta}^0 = (R^0,T^0,H^0,d^0)$ be the true value of $\underset{\sim}{\theta}$, and introduce

$$R_0 = R_0(\underset{\sim}{\theta}^0), \quad U_0 = R_1(\underset{\sim}{\theta}^0)/R_0, \quad V_0 = [U_0 + \tfrac{1}{2}],$$

$$h_0 = h_{V_0}(\underset{\sim}{\theta}^0),$$

$$M_0 = \min\Big\{x \in \{\alpha,\alpha+1, \cdots, \beta-1\}: h_x(\underset{\sim}{\theta}^0) \geq h_y(\underset{\sim}{\theta}^0)$$
$$\text{for all } y \in \{\alpha,\alpha+1, \cdots, \beta-1\}\Big\},$$

$$T_0 = R_0^2/\{\tfrac{4}{3}\pi(U_0-M_0)h_0^2\}, \quad d_0 = U_0 - T_0,$$

$$H_0 = h_0 T_0\sqrt{\pi}/R_0 \ .$$

Let $e$ be a $(\beta-\alpha) \times 1$-vector where all elements equal 1, let

$$\underset{\sim}{\iota} = R_0^{-1}(\alpha-T_0, \ \alpha+1-T_0, \ \cdots, \ \beta-1-T_0)',$$

and let

$$\underset{\sim}{\Phi} = (\underset{\sim}{e}, \ \underset{\sim}{e} \ T_0/R_0, \ \underset{\sim}{e} \ T_0 h_0 \sqrt{\pi}/R_0^2, \ \underset{\sim}{\psi}).$$

We then have

Theorem 8:  Under assumption 3,

$$N^{\frac{1}{2}}\{(\hat{R},\hat{T},\hat{H},\hat{d})' - (R_0,H_0,T_0,d_0)'\} \xrightarrow{\mathcal{L}_{\theta^0}} \mathfrak{N}(\underset{\sim}{0}, \underset{\sim}{\Phi}' \Sigma_0 \underset{\sim}{\Phi}).$$

Remark 1:  Note that $\underset{\sim}{\Phi}' \Sigma_0 \underset{\sim}{\Phi}$ is singular.

Remark 2:  Note also that $(R_0,T_0,H_0,d_0)$ is not the true value of the parameters here.  No one seems to have looked into the difference $(R^0,T^0,H^0,d^0) - (R_0,T_0,H_0,d_0)$ in any detail.

## 8.  A maximum likelihood method.

In sections $\underline{6}$ and $\underline{7}$, the estimators for $\underset{\sim}{\theta}$ appear as functions of the "raw" estimator $\hat{\underset{\sim}{\lambda}}$ for $\underset{\sim}{\lambda}$.  If one may really assume that $\underset{\sim}{\lambda}_a = \underset{\sim}{g}_a(\underset{\sim}{\theta}_a)$ for $a = 1,2, \cdots, A$, different approaches may be at least as efficient, however.  One obvious possibility is to enter the $\underset{\sim}{g}_a(\underset{\sim}{\theta}_a)$ into the likelihood function and maximize with respect to the $\underset{\sim}{\theta}_a$.  In the situation of section $\underline{3.4}$, this will amount to maximizing $\Sigma_a \hat{\eta}_a(\underset{\sim}{\theta}_a)$, where

$$\hat{\eta}_a(\underset{\sim}{\theta}_a) = \underset{d}{\Sigma} M(a,d) \ \ell n \ g_{ad}(\underset{\sim}{\theta}_a) - \underset{d}{\Sigma} L(a,d) \ g_{ad}(\underset{\sim}{\theta}_a) .$$

For simplicity, we shall assume that $\theta_{\sim 1}$, $\theta_{\sim 2}$, $\cdots$, $\theta_{\sim A}$ are functionally independent, so that we can maximize the likelihood function (if at all) by maximizing each $\hat{\eta}_a$ separately.

In the situation of section 4.3, the log likelihood function is of a different form, but we shall still construct an estimator $\theta_{\sim a}^*$ for $\theta_{\sim a}$ by maximizing $\hat{\eta}_a$.

The following theorem holds.

**Theorem 9:** Fix a $\epsilon\{1,2, \cdots, A\}$ and let the M(a,d) and L(a,d) be given as in section 3.4 or 4.3. Assume that $P\{L(a,d) > 0\} \to 1$ as $N \to \infty$ for all d where $(a,d)\epsilon_Q$. Make assumptions 1 and 2.

With a probability increasing to 1 as $N \to \infty$, there will then exist a value $\theta_a^{*(N)} \epsilon\Theta_a$ which maximizes $\hat{\eta}_a(\theta_{\sim a})$, and

$$N^{\frac{1}{2}}(\theta_{\sim a}^{*(N)} - \theta_{\sim a}^0) \xrightarrow{\quad \mathcal{L}\theta_{\sim a}^0 \quad} \mathcal{N}(\underset{\sim}{0}, \underset{\sim}{\Sigma}_{00}),$$

where $\underset{\sim}{\Sigma}_{00}$ is given by (6.4), provided there exist constants $k_a$ and $k_a'$ such that

$$(8.1) \quad k_a' \geq g_{ad}(\theta_{\sim a}) \geq k_a > 0 \text{ for all } \theta_{\sim a} \epsilon\Theta_a.$$

**Proof:** $1^0$. **Preliminaries.** Suppress the subscript a and fix the true value $\theta_{\sim}^0$. Let

$$\ell_d = \ell_d^{(N)} = E_{\theta_{\sim}^0} L(d), \quad L_d = \lim_{N\to\infty} \ell_d^{(N)}/N,$$

and note that [40, (10)]

$$E_{\underset{\sim}{\theta}^O} M(d) = \lambda_d^O \, \ell_d.$$

Let

$$\hat{Q}(\underset{\sim}{\lambda}) = \sum_d M(d) \, \ell n \, \lambda_d^* - \sum_d L(d) \, \lambda_d,$$

so that

$$\hat{\eta}(\underset{\sim}{\theta}) = \hat{Q}(g(\underset{\sim}{\theta})),$$

and let

$$Q(\underset{\sim}{\lambda}) = E_{\underset{\sim}{\theta}^O} \hat{Q}(\underset{\sim}{\lambda}) = \sum_d \ell_d \{\lambda_d^O \, \ell n \, \lambda_d - \lambda_d\},$$

and

$$\eta(\underset{\sim}{\theta}) = Q(g(\underset{\sim}{\theta})) = E_{\underset{\sim}{\theta}^O} \hat{\eta}(\underset{\sim}{\theta}).$$

Finally, let

$$\Delta(\underset{\sim}{\lambda}) = Q(\lambda^O) - Q(\underset{\sim}{\lambda}) = \sum_d \ell_d \{\lambda_d^O ( \ell n \, \lambda_d^O - \ell n \, \lambda_d) -$$
$$- (\lambda_d^O - \lambda_d)\}.$$

For large enough N, each $\ell_d$ will be positive. For such N, we will have $\Delta(\underset{\sim}{\lambda}) > 0$ for all $\underset{\sim}{\lambda} \neq \lambda^O$, and $\Delta(\underset{\sim}{\lambda})$ will strictly increase as each $|\lambda_d - \lambda_d^O|$ increases. For every $\epsilon > 0$ there then exists a $\delta_1(\epsilon)$ such that if $\Delta(g(\underset{\sim}{\theta})) \leqq \epsilon$ then $|g(\underset{\sim}{\theta}) - g(\underset{\sim}{\theta}^O)| \leqq \delta_1(\epsilon)$, and by the bicontinuity of $g$ there further exists a $\delta(\epsilon)$ such that $|\underset{\sim}{\theta} - \underset{\sim}{\theta}^O| \leqq \delta(\epsilon)$. Conversely there exists a $\delta_0(\epsilon)$ such that $\Delta(g(\underset{\sim}{\theta})) \leqq \delta_0(\epsilon)$ if $|\underset{\sim}{\theta} - \underset{\sim}{\theta}^O| \leqq \epsilon$. Let

$$S_\epsilon = \{\underset{\sim}{\theta}\epsilon\Theta: \; |\underset{\sim}{\theta}-\underset{\sim}{\theta}^0| \leq \epsilon\},$$

and choose $\epsilon$ so small that $S_\epsilon \subseteq \Theta$. Choose $\epsilon' > 0$, and let $\epsilon''$ be so small that $\epsilon'' \leq \delta_0(\epsilon)$, $\delta(2\epsilon'') \leq \epsilon'$, and

$$0 < 2\epsilon'' < \eta(\underset{\sim}{\theta}^0) - \inf \{\eta(\underset{\sim}{\theta}): \underset{\sim}{\theta}\epsilon\Theta\},$$

and let

$$\Theta_{\epsilon''} = \{\underset{\sim}{\theta}\epsilon\Theta: \; \Delta(g(\underset{\sim}{\theta})) \leq 2\epsilon''\}.$$

$2^0$. Existence of $\underset{\sim}{\theta}*$. Let

$$\gamma = \sup \Sigma_d \; \{|\ln g_d(\underset{\sim}{\theta})| + g_d(\underset{\sim}{\theta})\}.$$

By (8.1), $\gamma < \infty$. Let $A_{\epsilon''}^{(N)}$ be the event that

$$(8.2) \quad |N^{-1} M(d) - \lambda_d^0 \ell_d| < \epsilon'' / \gamma, \quad |N^{-1}L(d) - \ell_d| < \epsilon'' / \gamma.$$

Then $P_{\underset{\sim}{\theta}^0}(A_{\epsilon''}^{(N)}) \to 1$. Assume that (8.2) holds. Then

$$(8.3) \quad |\hat{\eta}(\underset{\sim}{\theta}) - \eta(\underset{\sim}{\theta})| < \epsilon'' \quad \text{for all } \underset{\sim}{\theta}\epsilon\Theta.$$

If $\underset{\sim}{\theta}\epsilon\Theta - \Theta_{\epsilon''}$, we therefore get

$$\hat{\eta}(\underset{\sim}{\theta}) < \eta(\underset{\sim}{\theta}) + \epsilon'' < \hat{\eta}(\underset{\sim}{\theta}^0) - \epsilon'' < \hat{\eta}(\underset{\sim}{\theta}^0),$$

so in maximizing $\hat{\eta}(\underset{\sim}{\theta})$ we need not take such $\underset{\sim}{\theta}$ into account. Since $\Theta_{\epsilon''}$ is closed and $\hat{\eta}(\cdot)$ is continuous, there exists a maximizing value $\underset{\sim}{\theta}* \epsilon \Theta_{\epsilon''}$.

$3^0$. Consistencey of $\underset{\sim}{\theta}*$. By the definition of

$\Theta_{\epsilon''}$, we have $\Delta(g(\theta^*)) \leq 2\epsilon''$. Thus $|\theta^*-\theta^0| \leq$

$\delta(2\epsilon'') \leq \epsilon'$, and the consistency of $\theta^*$ follows.

$4^0$. The theorem now follows from some general

results due to LeCam [50]. []

Remark: Note that $\Sigma_{00}$ is the most favourable asymptotic covariance matrix we can get by the procedures in section $\underline{6}$ when $\hat{\lambda}$ is given by (3.4). In this sense, therefore, the method of the present section is at least as good as any of the other general methods we have studied.

## 9. The choice of a graduating function.

$\underline{9.1.}$ In previous sections, it was presupposed that the applicability of a particular graduating function $g(\theta)$ had been established, and the problem was to estimate $\theta$. In many practical cases, the situation will be different. Instead of a single function $g$, there is often a finite family $\mathcal{F} = \{g(s,\theta) : s = 1,2, \cdots, S\}$ of candidates for a graduating function, and one is required to choose one of these on the basis of the data. In the case of human fertility rates, for instance, it is seldom given which function to use, and one may have to select one from among the Pearson family, the Hadwiger function, and the Brass polynomial (1.4), say.

We shall assume that all functions $g(s,\theta)$ have the same parameter space $\Theta$. This need not be the case originally, but it can be achieved by the introduction of dummy parameters if necessary.

9.2. To describe what it means to choose a function from the class $\mathfrak{F}$ "on the basis of the data," we shall assume that there is a member $g(s^0,\cdot)$ of $\mathfrak{F}$ which is the "true" graduating function. The choice of a member of $\mathfrak{F}$ then amounts to estimating $s^0$ as an extra parameter. A number of estimation procedures are in use (compare, e.g., [45, §6.5]),but their statistical properties do not seem to have been much investigated, except that one may know something about their consistency as $N \to \infty$. We list some of these procedures:

(i) For choosing among the members of the Pearson family, there exist standard methods [22], [57], [34], [13] based on the first four empirical moments. Keyfitz [45, p. 160] suggests that this type of criterion can also be used when the Hadwiger function (1.3) is included in $\mathfrak{F}$ along with the Pearson type functions.

(ii) In connection with the methods of section 6, an obvious procedure is to set

$$\hat{Q}_s(\theta) = N(\hat{\lambda} - g(s,\theta))' M(\hat{\lambda} - g(s,\theta)),$$

let $\hat{\theta}^{(s)}$ be a value of $\theta$ which minimizes $\hat{Q}_s(\theta)$

( for $s = 1, 2, \cdots, S$), and define $\hat{s}$ as the value of $\hat{s}$ that subsequently minimizes $\hat{Q}_s(\hat{\underset{\sim}{\theta}}^{(s)})$ [27], [71], [78], [79].

(iii) A similar criterion can be used in connection with the method of section $\underset{\sim}{8}$. Let

$$\hat{\eta}(s,\underset{\sim}{\theta}) = \sum_d \{M(d) \; \ell n \; g_d(s,\underset{\sim}{\theta}) - L(d) \; g_d(s,\underset{\sim}{\theta})\},$$

let $\underset{\sim}{\theta}^{*(s)}$ maximize this quantity, and let $\hat{s}$ be the value of $s$ that subsequently maximizes $\hat{\eta}(s,\underset{\sim}{\theta}^{*(s)})$.

(iv) Yntema [28], [79] has suggested calculating

$$\Delta_s = \sum_{d=1}^{D} |\hat{\lambda}_d - g_d(s,\hat{\underset{\sim}{\theta}}^{(s)})|$$

and

$$\Delta'_s = \max\{|\hat{\lambda}_d - g_d(s,\hat{\underset{\sim}{\theta}}^{(s)})| : \quad d = 1, 2, \cdots, D\},$$

and taking $\hat{s}$ as the $s$-value that minimizes $\Delta_s$ or $\Delta'_s$. Here $\hat{\underset{\sim}{\theta}}^{(s)}$ is any suitable estimator for $\underset{\sim}{\theta}$ based on $\underset{\sim}{g}(s,\underset{\sim}{\theta})$.

9.3. We shall take a look at the consistency properties of $\hat{s}$ as defined in 9.2(ii) and (iii) above.

By a proper specification of $\mathcal{I}$ and $\Theta$ we should be able to get $\underset{\sim}{g}(s',\Theta) \cap \underset{\sim}{g}(s'',\Theta)$ to be empty whenever $s' \neq s''$. (Otherwise, part of the values $\underset{\sim}{g}(s'',\theta)$ would be redundant.) To prove consistency, however, we need the stronger assumption that

(9.1)   $|g(s',\theta) - g(s'',\theta)| > 0$   for   $s' \neq s''$,

where   $|A-B| = \inf\{|a-b|: a \in A, b \in B\}$ denotes the Euclidean distance between two subsets   $A$   and $B$   of   $R_D$.

If (9.1) holds, if $\hat{\lambda}$ is consistent for $\lambda^0$, and if plim $M = M_0$ as $N \to \infty$, with $M_0$ positive definite, then $\hat{s}$ is consistent in section 9.2(ii) above.

Similarly, by step $1^0$ in the proof of theorem 9, $\hat{s}$ is consistent in section 9.2(iii) above when (9.1) holds.

9.4. Let $\{\hat{\theta}^{(s)}\}$ be some estimator which we would use for $\theta$ if it were known that $s^0 = s$, and assume that

(9.2)   $P_{s^0,\theta^0}\{N^{\frac{1}{2}}(\hat{\theta}^{(s^0)} - \theta^0) \in B\} \to \Phi(B)$ as $N \to \infty$,

where $\Phi$ is a limiting probability measure and $B$ is any $\Phi$-continuous measurable set. Let $\hat{s}$ be a consistent estimator for $s^0$. Then it is easy to show that

(9.3)   $P_{s^0,\theta^0}\{N^{\frac{1}{2}}(\hat{\theta}^{(\hat{s})} - \theta^0) \in B\} \to \Phi(B)$ as $N \to \infty$

for the same $B$. Of course, $\hat{\theta}^{(\hat{s})}$ is our estimator for $\theta$ and (9.3) tells us that its limiting distribution is the one we would get if $s^0$ were known. Similarly, $g(\hat{s}, \hat{\theta}^{(\hat{s})})$ will be our estimator for $\lambda^0$, and its asymptotic properties follow directly from (9.3).

## 10.   Concluding remarks.

10.1.   In the models described in sections 2.1 and 2.2 above, the seniority parameter is continuous. If it is known (or if one assumes) that one of the forces of transition can be represented by a nice and smooth parametric function, say

$$\lambda(x) = g(x; \underset{\sim}{\theta}),$$

and if one is faced with the problem of estimating $\lambda$, using analytic graduation is not necessarily the most obvious line of attack.  In fact, it seems more natural to try to construct an estimator $\hat{\underset{\sim}{\theta}}$ for $\underset{\sim}{\theta}$ directly, without going the way via the $\lambda^{\#}$, as described in section 3.  Grenander [29, pp. 76-91] has shown how this might be done for the force of mortality in example 2.2 (i) when the Gompertz-Makeham formula (1.1)(with continuous x) applies.  A similar investigation could be carried out for other forces of transition, like the forces of fertility of example 2.2(ii).

If one does not know enough about the function $\lambda(\cdot)$ to specify a parametric $g(\cdot, \theta)$ which can represent it, one may turn to nonparametric methods, such as those developed within reliability theory [5], [6].  The force of mortality in 2.2(i) appears there under the name of failure rate or hazard rate, and quite a lot of energy

has gone into finding suitable methods of estimating this function.

Although both of these types of approach were initiated by Grenander's paper [29] on mortality measurement, such techniques do not seem to be much in use in demography and related fields. One would be curious to know why this is so. Part of the explanation is, no doubt, that these developments are largely unknown among people working in those fields of application, but there are more valid reasons. We shall suggest some of them.

10.2. The following types of argument seem to be among the ones leading people to base their inferences from the data on the $M(a,d)$ and $L(a,d)$ only, and sometimes on the $\hat{\lambda}_{ad}$ only. (Note that least squares and moments method estimation procedures of $\underset{\sim}{\theta}$ only require knowledge of the $\hat{\lambda}_{ad}$.)

(i) In section 3.2 we described how the points $\{\zeta_d : d = 1,2, \cdots, D\}$ partitioning the seniority interval were selected according to conventional rules. Similarly, it is standard procedure to calculate "occurrence/exposure" rates of the kind developed in sections 3.4 and 4.3. The use of standard techniques, standard tabulations, and so on, facilitates comparison with other investigations of the same subject matter. This encourages the continued use of techniques which are already widely known and widely applied even when other methods may be known to a few people.

(ii) The reliability of the data which demographers have to work with, can be very weak due to phenomena such as age misreporting, underenumeration, and so on. Also, one frequently does not know more than approximate dates (e.g., the calendar year only) of occurrences of the events studies. This calls for the application of rather robust statistical techniques, such as those which we have described. Even though demographic data may be deficient, they may still be reliable enough to permit the use of the aggregated values $M(a,d)$ and $L(a,d)$, or at least the $\hat{\lambda}_{ad}$.

In many cases, the investigator does not even have access to the original data, but only to standard tabulations made from them. Such tabulations will often permit the use of methods described here, and rule out others.

(iii) A similar argument applies to the reliability of the models used. For example, most current models, including those considered in this paper, leave seasonal variation over the calendar year out of account. There is plenty of evidence of the importance of such variation in the occurrence of vital events, but in many cases this is just a nuisance factor which one wants to eliminate. Current methods relying on seniority interval lengths of at least a year seem to effectively do so.

(<u>iv</u>) Even in cases where the data are reliable and sufficiently detailed (and the present author believes that not nearly enough attention has been given to such cases), the information extracted by a statistical procedure should be geared to the needs of the user. It seems that a standard table of rates, like table 1, and certain other tables derived from it, contains just about as much information as can be handled in a substantive study. In fact, the prevalence of summary indices derived from such tables, and the extent to which argumentation is carried out in terms of such indices, suggests that the standard tables contain even too much information. The use of analytic graduation can be seen as another piece of evidence in the same direction, since it enables one to substitute the formula of a function and a (small) set of parameter values for a whole table. (This argument does not rule out the parametric procedure suggested at the beginning of section <u>10.1</u>.)

(v) Each of the estimators $\hat{s}$ listed in section <u>9.2</u> is a function of the data via $\hat{\underset{\sim}{\lambda}}$ only. This reflects the fact that an investigator faced with the problem of selecting a graduating function from a class $\mathfrak{F}$ of candidates is likely to calculate $\hat{\underset{\sim}{\lambda}}$, plot the corresponding diagram, and use this to decide which member to choose from $\mathfrak{F}$. In fact, this is the way in which certain graduating functions historically have been pinpointed as more suitable than others.

Once $\hat{s}$ has been determined, however, the investigator should not necessarily continue to use $\hat{\underset{\sim}{\lambda}}$ in the estimation of $\underset{\sim}{\theta}$, but should feel free to choose among all available procedures as far as the quality of his data permits.


10.3. It is probably appropriate to underline once more (compare section 1.1) that there exist many types of graduation methods in addition to analytic graduation techniques. Most of them were first developed for use in mortality studies, and in that context they are apparently applied at least as often as analytic methods are. Many of them must have been intended for use in other connections as well, for example in fertility studies. With the exception of graphic methods, however, their application to other types of vital rates than mortality rates seems much less popular. (Compare [53, p. 53].)

## ACKNOWLEDGEMENT

# References

[1]  H. Ammeter, "Wahrscheinlichkeitstheoretische
     Kriterien für die Beurteilung der Güte der
     Ausgleichung einer Sterbetafel," Mitt. Verein.
     schweizer. Versich.-Math., Vol. 52(1952),
     pp. 19-72.

[2]  H. Ammeter, "Der doppelseitige und die einseitigen
     $(I\chi^2)$-Tests und ihre Leistungsfähigkeit für die
     wahrscheinlichkeitstheoretische Überprüfung von
     Sterbetafeln," Blätter der Deutschen Gesellschaft
     für Versicherungsmathematik (Deutscher
     Aktuarverein)E. V., Vol. 1(1953), pp. 39-60.

[3]  T. W. Anderson, "An Introduction to Multivariate
     Statistical Analysis," New York, Wiley, 1958.

[4]  E. W. Barankin and F. Gurland, "On asymptotically
     normal, efficient estimators: I," University
     of California Publications in Statistics, Vol. 1
     (1951), pp. 89-130.  Compare [23].

[5]  R. E. Barlow, "Some recent developments in relia-
     bility theory," pp. 49-65 in Mathematisch Centrum
     Amsterdam, "European Meeting 1968, Selected
     Statistical Papers 2", Mathematical Center Tracts,
     No. 27, 1968.

[6]  R. E. Barlow and W. R. van Zwet, "Asymptotic

properties of isotonic estimators for the

generalized failure rate function. Parts I

and II," Reports No. ORC 69-5 and 69-10,

Berkeley, University of California, Operations

Research Center, 1969.


[7]  H. A. R. Barnett, "Graduation tests and experiments,"

J. Inst. Actuaries, Vol. 72(1951), pp. 15-54,

with discussion, pp. 55-74.


[8]  B. Benjamin, "Health and vital statistics," London,

George Allen & Unwin Ltd., 1968.


[9]  P. Bickel, "Lecture notes for asymptotic theory

(217B) Winter 1970," mimeographed.


[10]  P. Billingsley, "Convergence of Probability

Measures," New York, Wiley, 1968.


[11]  W. Brass, "The graduation of fertility distributions

by polynomial functions," Population Studies,

Vol. 14(1960), pp. 148-162.


[12]  J. M. Callies, "Utilisation de modèles mathéma-

tiques pour l'estimation des données démographiques

dans les pays en voie de developpement," Rev. Inst.

Internat. Statist.,Vol. 34 (1966), pp. 341-359.

[13] C. Chandrasekaran and P. P. Talwar, "Forms of age-specific birth rates by orders of birth in an Indian community," Eugenics Quarterly, Vol. 15 (1968), pp. 264-272.

[14] C. L. Chiang, "On regular best asymptotically normal estimates," Ann. Math. Statist.,Vol. 27(1956), pp. 336-351.

[15] C. L. Chiang, "Introduction to Stochastic Processes in Biostatistics,"New York, Wiley, 1968.

[16] H. Cramer and H. Wold, "Mortality variations in Sweden. A study in graduation and forecasting," Skand. Aktuarietidskr., Vol. 18(1935), pp. 161-240.

[17] M. Davies, "Linear approximation using the criterion of least total deviations," J. Roy. Statist. Soc. Ser. B, Vol. 29(1967), pp. 101-109.

[18] J. M. Dickey, "Smoothed estimates for multinomial cell probabilities," Ann. Math. Statist., Vol. 39 (1968) pp. 561-566.

[19] J. M. Dickey, "Smoothing by cheating," Ann. Math. Statist., Vol. 40(1969), pp. 1477-1482.

[20]   H. F. Dorn, "Methods of analysis for follow-up

studies," Human Biology, Vol. 22(1950),

pp. 238-248.


[21]   L. G. Du Pasquier, "Mathematische Theorie der

Invaliditätsversicherung," Mitt. Verein.

schweizer. Versich.-Math., Vol. 7(1912),

pp. 1-7, and Vol. 8(1913), pp. 1-153.


[22]   W. P. Elderton and N. L. Johnson, "Systems of

Frequency Curves," London, Cambridge University

Press, 1969.


[23]   G. U. Fenstad, "A note on a theorem of Barankin

and Gurland," University of Oslo, Institute

of Mathematics, Statistical Research Report

No. 2, 1965.


[24]   T. S. Ferguson, "A method of generating best

asymptotically normal estimates with application

to the estimation of bacterial densities,"Ann.

Math. Statist.,Vol. 29(1953), pp. 1046-1062.


[25]   E. Fix and J. Neyman, "A simple stochastic model

of recovery, relapse, death and loss of patients,"

Human Biology, Vol. 23(1951), pp. 205-241.

[26] D. A. S. Fraser, "Nonparametric Methods in
        Statistics," New York, Wiley, 1957.

[27] E. Gilje, "Fitting curves to age-specific fertility
        rates. Some examples," Statistical Review,
        Ser. III, Vol. 7(1969), pp. 118-134, Swedish
        Central Bureau of Statistics, Stockholm.

[28] E. Gilje and L. Yntema, "The shifted Hadwiger
        fertility function," to appear.

[29] U. Grenander, "On the theory of mortality
        measurement," Skand. Aktuarietidskr., Vol. 39
        (1956), pp. 70-96 and 125-153.

[30] R. D. Grove and A. M. Hetzel, "Vital statistics
        rates in the United States, 1940-1960,"
        National Center for Health Statistics, U.S.
        Department of Health, Education and Welfare,
        1968.

[31] H. Hadwiger, "Eine analytische Reproduktions-
        funktion für biologische Gesamtheiten,"
        Skand. Aktuarietidskr., Vol. 23(1940),
        pp. 101-113.

[32] T. E. Harris, P. Meier, and J. W. Tukey, "Timing
        of the distribution of events between observa-

tions," <u>Human Biology</u>, Vol. 22(1950),

pp. 249-270.


[33]  <u>L. Henry</u>,          "D'un problème fondamental

de l'analyse démographique," <u>Population</u>,

Vol. 14(1959), pp. 9-32.


[34]  <u>A. B. Hoadley</u>, "Use of the Pearson densities for

approximating a skew density whose left

terminal and first three moments are known,"

<u>Biometrika</u>, Vol. 55(1968), pp. 559-563.


[35]  <u>J. M. Hoem</u>, "Markov chain models in life insurance,"

<u>Blätter der Deutschen Gesellschaft für</u>

<u>Versicherungsmathematik (Deutscher Aktuarverein)</u>

<u>E. V.</u>, Vol. 9(1969), pp. 91-107.


[36]  <u>J. M. Hoem</u>, "Fertility rates and reproduction

rates in a probabilistic setting," <u>Biométrie-</u>

<u>Praximétrie</u>, Vol. 10(1969), pp. 38-66.


[37]  <u>J. M. Hoem</u>, "A probabilistic model for primary

marital fertility," <u>Yearbook of Population</u>

<u>Research in Finland</u>, Vol. 11(1969), pp. 73-86.


[38]  <u>J. M. Hoem</u>, "Probabilistic fertility models of the

life table type," <u>Theoretical Population</u>

<u>Biology</u>, Vol. 1(1970), pp. 12-38.

[39]  J. M. Hoem, "Some results concerning the number
      of jumps in a finite time interval of a
      Markov chain with a continuous time parameter,"
      University of Oslo, Institute of Economics,
      Memorandum of February 2, 1970.

[40]  J. M. Hoem, "Point estimation of forces of transi-
      tion in demographic models," J. Roy. Statist.
      Soc. Ser. B, Vol. 32(1970), to appear.

[41]  J. M. Hoem, "A probabilistic approach to nuptiality,"
      Biométrie-Praximétrie, to appear.

[42]  L. Høyland, "Estimation in follow-up studies,"
      University of Oslo, Institute of Mathematics,
      Statistical Research Report No. 4, 1967.

[43]  H. Jecklin and P. Strickler, "Wahrscheinlich-
      keitstheoretische Begründung mechanischer
      Ausgleichung und deren praktische Anwendung,"
      Mitt. Verein. Schweizer. Versich.-Math.,
      Vol. 54(1954), pp. 125-161.

[44]  D. A. Jones, "Bayesian statistics," Trans. Soc.
      Actuaries, Vol. 17(1965), pp. 33-57.

[45]  N. Keyfitz, "Introduction to the Mathematics of
      Population," Reading, Mass., Addison-Wesley,
      1968.

[46]  G. S. Kimeldorf and D. A. Jones, "Bayesian
          graduation," Trans. Soc. Actuaries, Vol. 19
          (1967), pp. 66-112.

[46a] Yu.L. Kruopis, "On estimates of transition
          intensities with migration," Theory Probab.
          Appl.,Vol. 14(1969), pp. 219-228.

[47]  I. Lah, "Generalization of Yastremsky's formula
          for analytical graduation of fertility rates,"
          J. Roy. Statist. Soc. Ser. A, Vol. 121(1958),
          pp. 100-104.

[48]  I. Lah, "Analytische Ausgleichung der aus den
          Ergebnissen der Volkszählungen berechneten
          demographischen Tafeln," International Union
          for the Scientific Study of Population, Con-
          ference, Wien, 1959, pp. 192-201.

[49]  I. Lah, "Analytically graduated fertility of married
          women in Australia with respect to the duration
          of marriage," pp. 266-276 in International Union
          for the Scientific Study of Population,
          "Contributed Papers, Sidney Conference, Australia,
          21st to 25th August, 1967."

[50]  L. LeCam, "Asymptotic least squares theory,"
          mimeographed.

[51]  S. Lipton and C. A. McGilchrist, "The derivation

of methods for fitting exponential regression

curves," Biometrika, Vol. 51(1964), pp. 504-507.

[52]  A. J. Lotka, "Theorie analytique des associations

biologiques. Part II. Analyse démographique

avec application particulière à l' espèce

humaine," Paris, Hermann & Cie, 1939.

[53]  D. P. Mazur, "The graduation of age-specific

fertility rates by order of birth of child,"

Human Biology, Vol. 39(1967), pp. 53-64.

[54]  M. D. Miller (et al.), "Elements of Graduation,"

New York, Actuarial Society of America, 1942.

[55]  S. Mitra, "The pattern of age-specific fertility

rates," Demography, Vol. 4(1967), pp. 894-906.

[56]  H. D. Patterson, "A further note on a simple method

for fitting an exponential curve," Biometrika,

Vol. 47(1960), pp. 177-180.

[57]  E. S. Pearson and H. O. Hartley (Eds.), "Biometrika

Tables for Statisticians," London, Cambridge

University Press, 1966.

[58]  W. Saxer, "Versicherungsmathematik. Zweiter Teil,"

Berlin, Springer-Verlag, 1958.

[59]   B. Schneider. "Die Bestimmung der Parameter im
       Ertragsgesetz von E. A. Mitscherlich," Biom.
       Zeit.,Vol. 5(1963), pp. 78-95.

[60]   H. L. Seal, "Tests of a mortality table graduation,"
       J. Inst. Actuaries, Vol. 71(1941), pp. 5-47,
       with discussion, pp. 48-67.

[61]   M. C. Sheps, "Characteristics of a ratio used to
       estimate failure rates:  occurrences per person
       year of exposure," Biometrics, Vol. 22(1966),
       pp. 310-321.

[62]   M. C. Sheps, "On the person year concept in
       epidemiology and demography," Milbank Memorial
       Fund Quarterly, Vol. 44(1966), pp. 69-91.

[63]   M. C. Sheps, J. A. Menken, and A. P. Radick,
       "Probability models for family building, an
       analytical review," Demography, Vol. 6(1969),
       pp. 161-183.

[64]   M. C. Sheps and J. A. Menken, "On closed and open
       birth intervals in a stable population,"
       paper prepared for the Segunda Conferencia
       Regional de Poblacion, Mexico City, August, 1970.

[65]  W. Simonsen, "Forsikringsmatematik. Hefte I og II,"
      Köbenhavns Universitets Fond til Tilvejebringelse
      af Laeremidler, 1966-67.

[66]  M. Spiegelman, "Introduction to Demography.
      Revised Edition," Cambridge, Mass., Harvard
      University Press, 1968.

[67]  W. L. Stevens, "Asymptotic regression," Biometrics,
      Vol. 7(1951), pp. 247-267.

[68]  E. Sverdrup, "Basic concepts in life assurance
      mathematics," Skand. Aktuarietidskr., Vol. 35
      (1952), pp. 115-131.

[69]  E. Sverdrup, "Estimates and test procedures in
      connection with stochastic models for deaths,
      recoveries and transfers between different
      states of health," Skand. Aktuarietidskr.,
      Vol. 46(1965), pp. 184-211.

[70]  E. Sverdrup, "Laws and Chance Variations. Vol. 1,"
      Amsterdam, North-Holland, 1967.

[70a] P. P. Talwar, "Age patterns of fertility,"
      University of North Carolina at Chapel Hill,
      Institute of Statistics Mimeo Series No. 656,
      1970.

[71] K. Tekse, "On demographic models of age-specific fertility rates," Statistical Review, Ser. III, Vol. 5(1967), pp. 189-207, Swedish Central Bureau of Statistics, Stockholm.

[72] E. J. Wegman, "Maximum likelihood histograms," University of North Carolina at Chapel Hill, Institute of Statistics Mimeo Series No. 629,1969.

[73] E. J. Wegman, "Nonparametric probability density estimation," University of North Carolina at Chapel Hill, Institute of Statistics Mimeo Series No. 638, 1969.

[74] K. Weichselberger,"Über eine Theorie der gleitende Durchschnitte und verschiedene Anwendungen dieser Theorie," Metrika, Vol. 8 (1964), pp. 185-230.

[75] S. D. Wicksell, "Nuptiality, fertility, and reproductivity," Skand. Aktuarietidskr., Vol. 14(1931), pp. 125-157.

[76] R. A. Wijsman, "On the theory of BAN - estimates," Ann. Math. Statist., Vol. 30(1959), pp. 185-191, correction note, pp. 1268-1270.

[77]  H. H. Wolfenden, "The fundamental Principles of
      Mathematical Statistics (with special reference
      to the requirements of actuaries and vital
      statisticians, and an outline of a course in
      graduation)," New York, Actuarial Society of
      America, 1942.

[78]  L. Yntema, "The graduation of net fertility tables,"
      Boletim do Instituto dos Actuários Portugueses,
      Vol. 8(1953), pp. 29-43.

[79]  L. Yntema, "On Hadwiger's fertility function,"
      Statistical Review, Ser. III, Vol. 7(1969),
      pp. 113-117, Swedish Central Bureau of
      Statistics, Stockholm.

[80]  S. Zahl, "A Markov process model for follow-up
      studies," Human Biology, Vol. 27(1955),
      pp. 90-120.

[81]  E. Zwinggi, "Versicherungsmathematik," Basel and
      Stuttgart, Birkhäuser Verlag, 1958.

Figure 1.

Age-specific mortality rates per 1000.
Females, Oslo, 196 .

Observed rates:
$1000y = 0.113 + 0.0079(1.12074)^X$:
Gompertz - Makeham function, fitted
by minimum $\chi^2$

Figure 2.

Age-specific fertility rates per 1000.

Females, Stavanger (Norway), 1963.