

# A SYSTEM OF SOCIAL MATRICES

BY RICHARD STONE

*Cambridge, England*

The paper is concerned with a method of organizing and analyzing information relating to human stocks and flows. The kind of statistical reporting system envisaged is of a traditional kind, but extended so as to record year-to-year changes of state. Life is divided into a number of sequences, each with its own set of characteristic classifications, to avoid an excessive proliferation of categories and so enable many analyses to be made with the kind of statistics already available in a number of countries. The need, for some analytical purposes, to combine classifications from different sequences is fully recognized; and this need indicates a direction in which statistical reporting systems should move in the future.

The main analytical tool is a set of linear difference equations which, under suitable conditions, can be interpreted either in terms of an input-output system, as in economics, or in terms of an absorbing Markov chain, as in probability theory. A simple regression model is used to link characteristic classifications.

About half the paper is taken up with numerical examples, mainly connected with the British educational system as it was in the mid-1960's. An application is also given to movements into, through and out of a psychiatric service system in Scotland.

## I. INTRODUCTION

Since writing the paper which I contributed four years ago to the meeting of the International Association for Research in Income and Wealth at Maynooth [11], I have continued to work in the area of social demography. Part of what I have written in this period has been connected with the endeavour promoted by the Statistical Office of the United Nations to formulate an integrated system of demographic, manpower and social statistics. In early discussions on the subject it was thought that, whatever else such a system should include, it must contain three ingredients. The first of these is a detailed treatment of human stocks and flows in different areas of social interest, such as education, employment, health, delinquency and so on. The second is a means of accounting for the services provided, often by the state, in these areas, the costs incurred and the resources engaged in them. And the third is a means of recording the distribution of these services over various classes of beneficiary. A fairly lengthy discussion paper in which an attempt was made to elaborate and illustrate these ideas is available in [12].

In this paper I shall concentrate on the socio-demographic aspect of this system, on the development of what may be called life sequences; but I shall also indicate how this information can be linked with economic information on costs and benefits. I shall give some numerical examples drawn from different fields and try to answer a number of general questions to which the approach gives rise.

## II. CHARACTERISTICS AND STATES

In his progression from birth to death an individual passes through a succession of states. Each year he becomes a year older; at some time between the

ages of two and five he is almost certain to start going to school; at any time his parents may move to another part of the country; at any time he may fall ill and require to be treated by a doctor or to spend time in a hospital; at any time his aberrant behaviour may turn to delinquency and, beyond a certain age, this will be recognized by the criminal courts and, if detected, bring him in contact with the penal system.

Apart from these characteristics, of which some, like age, must change with time in a perfectly regular way, some, like health and educational attainment, are bound to change with time but not in such a regular way, and some, like social behaviour, may or may not change with time, there are other characteristics either of the individual or of his family which can hardly change. The individual is male or female, black or white, clever or stupid, tidy or untidy, handsome or ugly, to mention but a few personal characteristics. His family is patrician or plebeian, rich or poor, religious or irreligious, strict or easy going, settled or migrant, to mention but a few family characteristics.

Obviously an individual can be described by a very large number of characteristics. Any attempt to classify individuals by many characteristics simultaneously leads, as the number increases, to more and more compound categories, that is states, and calls for large quantities of data. While it must be recognized that some analytical purpose may require a classification by almost any combination of characteristics, a regular statistical reporting system must almost inevitably provide information on a less ambitious scale. For any aspect of life, information can be confined to what is thought to be necessary to describe that aspect, and the merging of information relating to different aspects can be treated as a separate problem. This application of the principle of divide and conquer leads naturally to the concept of life sequences.

### III. LIFE SEQUENCES

A life sequence traces the changes of state from birth to death in some particular compartment of life. For purposes of statistical reporting this requires that we define the compartment in question and that we draw up a list of the classifications to be considered characteristic of it.

In defining the compartment of life to which a sequence relates it is convenient to make use of the concept of a boundary as exemplified by the economic concept of the production boundary. Applying this idea to education, for instance, we might decide to draw the boundary round full-time formal education, say, and ignore all types of part-time and informal education; or we might decide to extend the boundary so as to include some but not all of these peripheral types of education. The need for this choice and the way it is made finds a close parallel in the definition of production. In the present instance we could defend the tight definition of the educational boundary by reference to the useful results that can be obtained from a study of full-time formal education and to the difficulties of collecting information about part-time and informal education.

In formulating the characteristic classifications of a sequence, the main endeavour should be to provide an adequate description of what takes place in it. For many purposes, age and sex should be considered characteristic classifi-

cations in all sequences, although this is not essential. In the case of learning activities, additional classifications which would obviously be desirable are: type of educational institution attended, grade or level of work, subjects studied and leaving qualifications.

Every sequence relates to the whole population of a country or region, whether the data are provided by complete enumeration or by sampling. Typically, therefore, there will always be one or more "inactive" categories. For instance in the learning sequence, it would be necessary to record those who were educationally inactive and to divide those who had not yet entered the educational system from those who had left it.

Let us now turn to some examples of life sequences.

(a) *The active sequence.* This sequence relates to learning activities, earning activities and the educationally and economically inactive. Its various parts are discussed at some length in [9]. A possible boundary and set of classifications for learning activities have been given above; for earning activities, the boundary would be drawn by reference to the concept of production in the SNA and obvious characteristic classifications are occupation, industrial status and industry of employment. The remainder of the population comes into the category "inactive".

(b) *The passive sequence.* This sequence relates to the succession of family groupings to which individuals are attached in the course of their life. The great majority of individuals are attached to natural households of different kinds but a sizeable minority are attached at some period of their life to one or other of the many forms of institutional household. Characteristic classifications for natural households are, for example; size, composition, social class, income, race, religion, housing conditions, neighbourhood and location.

(c) *The sequence of health and medical care.* This sequence relates to health conditions, their treatment and their consequences. Such information may be useful in two different contexts: (i) in planning and organizing health services; and (ii) in studying the aetiology and treatment of diseases. In the first case, the boundary will be drawn with reference to the organizational complex in question which may stretch from a national health service to a group of hospitals, clinics, and practitioners engaged in providing a particular kind of health service in a particular locality. In the second case, the boundary will be drawn with reference to a disease or group of diseases. Classifications characteristic of this sequence are conditions of health, diseases, treatments, incapacity and medical practitioners and establishments.

(d) *The sequence of delinquency.* This sequence relates to aberrant behaviour leading to delinquency, to delinquency itself and to its treatment and consequences. As in the case of health and medical care, we can interest ourselves in the organizational problems of dealing with delinquency in general or in the aetiology and treatment of particular groups of crimes; and the boundary must be drawn by reference to these interests. Classifications characteristic of this sequence are early warning behaviour, offences, gravity of offences, treatment of offenders, incapacity and institutions (police, courts, prisons, etc.) which deal with offences and offenders.

Eventually we shall consider the problem of connecting classifications from different sequences. But before we do this let us look at the questions of presenting and analyzing data relating to a single sequence. For many purposes this will be sufficient; for instance, in planning educational or health services we are probably more interested in the number of individuals likely to move through different parts of the system in the future than with, say, the personal or familial characteristics of these individuals. These characteristics are mainly relevant for another purpose, namely in understanding why some kinds of individual tend to be concentrated in certain parts of the system. Of course, this knowledge may help us to make better projections but it is not necessarily essential for that purpose.

#### IV. A FRAMEWORK FOR HUMAN STOCKS AND FLOWS

In studying any sequence we need to know how individuals are distributed over states at different points in time (stock information) and how individuals move between states over intervals of time (flow information). Contrary to the position in economics, the world of social statistics contains a great deal of stock information and comparatively little flow information. While stock information has many uses it does not enable us to see in any detail how changes take place. For this reason it seems to me that the improvement of flow information is the most important task, at present, in this corner of the statistical universe. While I think it useful, by formulating models, to show how statistical information can be used, I also think that a coherent measurement of stocks and flows is elementary, and desirable independently of any particular system of models, such as the one described in the next section.

A standard matrix framework which can accommodate the data for any sequence is set out in table 1 below. This framework is a little different from the one I used in my earlier paper [11] since it relates to the opening and closing stocks of one period rather than to the outflow from two successive periods. This distinction is spelled out in [9, 12].

This framework can accommodate information from any sequence or indeed, from any combination of sequences; everything depends on the classifications used. The information itself can relate to those alive in a given interval of time (cross-section or transversal data) or to those born in a given interval of time (vintage or longitudinal data).

There are several advantages in having a framework of general applicability. In the first place, it enables us to specify sets of data which are coherent and, as far as they go, complete. In the second place, this feature is of particular importance if, as frequently happens, the data for a single sequence are collected by separate agencies. For instance, it enables us to see precisely what demographic, educational, employment and other statistics must be available to construct a coherent matrix for the active sequence. In the third place, the eventual crossing of classifications from different sequences is likely to be made easier if any common parts of different sequences are compatible. In the fourth place, a framework is useful in building models because it shows the identities by which the variables in the model are connected and so the degrees of freedom available

TABLE 1  
THE STANDARD MATRIX

State at New Year $\theta + 1$ \ State at New Year $\theta$	Outside world	Our country: opening states	Closing stocks
Outside world	$\alpha$	$d'$	
Our country: closing states	$b$	$S$	$\Lambda n$
Opening stocks		$n'$	

The symbols in this table have the following meaning:

- $\alpha$ , a scalar (or single number), denotes the total number of individuals who both enter and leave "our country" in the course of the period and so are not recorded in either the opening or the closing stock. An example is a baby born during the period who dies before the end of the period.
- $b$ , a column vector (or column of numbers), denotes the new entrants into "our country", namely the births and immigrations of the period, who survive to the end of the period. Individuals in this category are recorded in the closing stock but not in the opening stock.
- $d'$ , a row vector (or row of numbers), denotes the leavers from "our country", namely the deaths and emigrations of the period. Individuals in this category appear in the opening stock but not in the closing stock. In accordance with convention, a row vector is represented by the symbol for a column vector followed by a prime superscript (').
- $S$ , a square matrix (or square block of numbers), denotes the survivors in "our country" through the period, who are recorded in both the opening and the closing stock. They are classified by their opening states in the columns and by their closing states in the rows.
- $n'$ , a row vector, denotes the opening stock in each state.
- $\Lambda n$ , a column vector, denotes the closing stock in each state. The symbol  $\Lambda$  denotes the lag operator which shifts in time the variable to which it is applied. Thus, if  $n(\tau)$  denotes the value of the vector  $n$  at time  $\tau$ , then  $\Lambda n(\tau) \equiv n(\tau + 1)$  and, in general  $\Lambda^\theta n(\tau) \equiv n(\tau + \theta)$ .

to be absorbed by behavioural relationships, policy constraints and the like. Finally, even if our statistics came from a continuously updated system of individualized data, it would still be necessary to have clear ideas about the information to be extracted from the data bank for any particular analysis and the identities connecting these data. Of course, in these circumstances, the arrangement of data on worksheets or printed pages ceases to be of interest; they are all in the computer. But the processes of thought needed to use them effectively are much the same as when simpler methods of processing are used.

## V. SOME SIMPLE MODELS

The symbols in Table 1 above are connected by two sets of identities.

First, those in the opening stock either survive to the end of the year or die or emigrate in the year, that is

$$(1) \quad n = S'i + d$$

where  $i$  denotes the unit vector, so that  $S'i$  denotes the column sums of  $S$ . Second, those in the closing stock either have survived from the beginning of the year or are born or immigrate during the year, that is

$$(2) \quad \Lambda n = Si + b.$$

(a) *Backward models and forward models.* Either of these equations can be turned into a set of different equations with an exogenous vector by forming coefficient matrices based respectively on the rows or columns of  $S$ . If the row coefficient matrix is denoted by  $G'$ , then

$$(3) \quad G' = S'\Lambda\hat{n}^{-1}$$

and so, by combining (1) and (3),

$$(4) \quad n = G'\Lambda n + d.$$

Similarly, if the column coefficient matrix is denoted by  $C$ , then

$$(5) \quad C = S\hat{n}^{-1}$$

and so, by combining (2) and (5),

$$(6) \quad \Lambda n = Cn + b$$

(b) *The forward model in conditions of stationary equilibrium.* Equation (6) is based on transition proportions: those in a given state at the beginning of a year go in fixed proportions to the states that it directly feeds. This equation enables us to project forward. In what follows I shall concentrate mainly on this form of equation; *mutatis mutandis* parallel statements can always be made for the backward equation.

Let us consider first the case in which our data relate to a population in stationary equilibrium or can in some way be adjusted to satisfy this condition. In these circumstances the size and composition of the population remains unchanged and so  $\Lambda^{\theta}n = n$  and  $\Lambda^{\theta}b = b$ . In this case (6) takes the form.

$$(7) \quad \begin{aligned} n &= Cn + b \\ &= (I - C)^{-1}b \end{aligned}$$

which expresses  $n$  as a matrix transform of  $b$ .

To an economist, (7) is formally identical to the quantity equation of an open input-output system, with  $(I - C)^{-1}$  as the matrix multiplier that transforms final demands into total outputs. In the present case, the matrix transforms new entrants into total population. The significance of this interpretation is that, as in the economic case, there is a corresponding price equation which enables us to work out future costs or revenues associated with those now in any given state from a knowledge of the unit cost or revenue in each state.

To a probability theorist, the matrix inverse  $(I - C)^{-1}$ , has the form of the fundamental matrix of an absorbing Markov chain; and it can be given this interpretation if  $C$  can be regarded as a probability matrix and not merely as a

matrix of proportions. This implies that the probability of movement from a given state to another state is the same for all members of the given state, which in turn implies that the probabilities of movements from the given state are independent of the path by which that state has been reached. The significance of this interpretation is that if  $C$  can be regarded as a probability matrix then the sequence can be regarded as an absorbing Markov chain and the many theorems applicable to such chains can be applied to it.

(c) *Remembering the past.* If states are defined in terms of the current characteristics of individuals only, it may not be plausible to assume that the probabilities of movements from states are independent of the paths along which those states have been reached. For instance, in the sequence of health and medical care, we may find two individuals in apparently the same conditions of health at a given age but we might expect their medical futures to be different if their medical pasts had been different. In such a case it would seem necessary so to define states that they relate not only to the medical situation of the moment but also to the past medical situations of the individuals. This can be done as follows.

Suppose that the life span is divided into  $\tau$  age groups or stages, and that at each stage individuals are classified to  $\mu$  medical categories. In the simplest case the medical categories might consist of the dichotomy well or ill. There would in this case be two states at the first stage. At the second stage, those who were well at the first stage would be classified according as they were well or ill and those who were ill at the first stage would be similarly classified. Thus, at the second stage there would be  $\mu^2$  states and, in general, at stage  $\tau$  there would be  $\mu^\tau$  states.

If we think in terms of a period of one year between the opening and closing stocks and of a stage length of ten years, then in any period an individual can: (i) remain in the stage and, as far as this paper is concerned, also the state in which he was recorded at the beginning of the period; (ii) move to one of the states characteristic of the next stage; or (iii) move into the absorbing state.

With this method of recording, the  $C$ -matrix takes a very special form: the diagonal submatrices are diagonal (corresponding to the fact that changes of state within a stage are not recorded); and the only other non-zero submatrices are those immediately below the diagonal ones (corresponding to the fact that individuals can only go from one stage to the next and can neither skip a stage nor go backwards). Thus, in the case of three stages,  $C$  takes the form

$$(8) \quad C = \begin{pmatrix} \hat{c}_{11} & 0 & 0 \\ C_{21} & \hat{c}_{22} & 0 \\ 0 & C_{32} & \hat{c}_{33} \end{pmatrix}$$

whence

$$(9) \quad (I - C)^{-1} = \begin{bmatrix} (I - \hat{c}_{11})^{-1} & 0 & 0 \\ (I - \hat{c}_{22})^{-1}C_{21}(I - \hat{c}_{11})^{-1} & (I - \hat{c}_{22})^{-1} & 0 \\ (I - \hat{c}_{33})^{-1}C_{32}(I - \hat{c}_{22})^{-1} & (I - \hat{c}_{33})^{-1}C_{32}(I - \hat{c}_{22})^{-1} & (I - \hat{c}_{33})^{-1} \\ \quad \times C_{21}(I - \hat{c}_{11})^{-1} & & \end{bmatrix}$$

Thus, while the  $C$ -matrices tend to be large, the inverse matrices can be built up by taking reciprocals and by systematic matrix multiplication.

(d) *The forward model as a basis for projections.* If it can be assumed that the  $C$ -matrix remains fixed over time, then (6) can be used to make projections provided that we know the future course of the exogenous vector,  $b$ . Thus, if we apply the lag operator,  $\Lambda$ , to (6), we obtain

$$(10) \quad \begin{aligned} \Lambda^2 n &= C\Lambda n + \Lambda b \\ &= C^2 n + Cb + \Lambda b \end{aligned}$$

and, in general,

$$(11) \quad \Lambda^\tau n = C^\tau n + \sum_{\theta=0}^{\tau-1} C^\theta \Lambda^{\tau-\theta-1} b$$

Equation (11) expresses the stock vector  $\tau$  periods hence in terms of the present stock vector and new entry vectors from the present period through period  $\tau-1$ .

(e) *Changing coefficients.* The elements of the  $C$ -matrix, like those of the  $A$ -matrix in economic input-output analysis, depend on supplies and demands which, in turn, depend largely on public policy and on the community's attitude to education, health or whatever it may be. It is likely, therefore, that the  $C$ -matrix will change over time, and the question arises: can we find a satisfactory method of projecting the elements of  $C$  so that, as we move forward, we can gradually change the transition probabilities to be applied to opening stock vectors? If we can then, as is shown in [9, 12], there is no difficulty in reformulating (11) to incorporate this information.

If we look at a series of  $C$ -matrices we find that apart from sudden changes, occasioned by such policy decisions as raising the school-leaving age, the transition probabilities are either constant or changing slowly. For instance, the probabilities of remaining at school at ages following the school-leaving age are rising and, since those who leave tend more and more to go on to some form of further education, the probability of seeking employment at these ages is falling. However, even if we can measure these probabilities over the last twenty years, we do not possess a very secure base for projecting over the next twenty years and so we can only use simple methods which amount to little more than trend projections with allowance for expected sudden changes. The method with which I have experimented is the simple epidemic model applied to educational transitions, the transition at any age to employment being treated as a residual. If  $c_{sr}$  denotes the transition probability from educational state  $r$  to educational state  $s$ , this model, when expressed in terms of discrete time, takes the form

$$(12) \quad \begin{aligned} \Delta c_{sr} &= \beta c_{sr} (\gamma - c_{sr}) \\ &= \beta \gamma c_{sr} - \beta c_{sr}^2 \end{aligned}$$

where  $\Delta \equiv \Lambda - 1$  and  $0 < \gamma \leq 1$  denotes the maximum value that  $c_{sr}$  can take. If (12) is written in continuous time, its integral is a logistic curve and so, with time,  $c_{sr}$  will move towards  $\gamma$  and cannot take on impossible values as it could if it were assumed to move along a linear or an exponential time trend. It will



often be found that the data are insufficient to determine  $\gamma$  at all accurately and in this case it will be necessary to select arbitrary values of  $\gamma$  out of a plausible range and see how much difference the selection makes to projections twenty or thirty years into the future.

(f) *The price equation in the forward model.* I propose to discuss this equation in terms of educational costs though it can, of course, be applied to any other costs or, indeed, to gains or the excess of gains over costs of any kind whatsoever.

Let  $m$  denote a vector whose elements measure the educational costs that must be incurred this year to educate an individual now in a given state of the system. On the assumption that  $m$  remains fixed in the future, the total cost to be incurred from now on to educate, or complete the education of, an individual now in a given state is an element of a vector,  $k$  say, where

$$(13) \quad \begin{aligned} k &= m + C'm + C'^2m + \dots \\ &= m + C'k \\ &= (I - C')^{-1}m \end{aligned}$$

The terms on the right-hand side of the first row of (13) relate to the successive years in which educational costs will be incurred. The elements of these vectors relate to the present states of individuals multiplied by the probable educational costs they will incur this year, next year and so on.

If it can be assumed that unit costs will change so that, in year  $\theta$ ,  $m$  will be replaced by  $\Lambda^\theta m$ , then (13) becomes

$$(14) \quad \begin{aligned} k &= m + C'\Lambda m + C'^2\Lambda^2m + \dots \\ &= m + C'\Lambda k \\ &= (I - C'\Lambda)^{-1}m \end{aligned}$$

Thus if we can estimate  $\Lambda^\theta m$  for the relevant values of  $\theta$  we can allow for changing costs.

If it is assumed that the  $C$ -matrix will change,  $C^\theta$  must be replaced by  $\bar{C}^\theta \equiv \Lambda^{\theta-1}C \cdot \Lambda^{\theta-2}C \dots C$ . If we define

$$(15) \quad \bar{C} \equiv I - \left[ I + \sum_{\theta=1}^{\infty} \bar{C}^\theta \right]^{-1}$$

then we can write

$$(16) \quad \begin{aligned} k &= m + C'\Lambda m + (\Lambda C \cdot C')\Lambda^2m + \dots \\ &= (I + \bar{C}'\Lambda + \bar{C}'^2\Lambda^2 + \dots)m \\ &= (I - \bar{C}\Lambda)^{-1}m \end{aligned}$$

Thus if we can estimate  $\Lambda^\theta C$  for the relevant values of  $\theta$  we can allow for changing transition probabilities.

If  $\rho$  denotes the rate of interest,  $\sigma \equiv 1/(1 + \rho)$  denotes the discount factor and if the states of  $C$  are separated by annual intervals (as would be the case if year of birth were the primary criterion of classification) it is easy to calculate the discounted streams of future costs corresponding to (13). If  $k^*$  denotes the vector of discounted accumulated costs and if  $C^* \equiv \sigma C$ , then (13) is replaced by

$$(17) \quad k^* = (I - C^*)^{-1}m$$

If we have calculated the inverse in (13), we can readily calculate the inverse in (17) since

$$(18) \quad (I - C^*)^{-1} \equiv \hat{s}(I - C)^{-1}\hat{s}^{-1}$$

where the elements of  $s$  are descending powers of  $\sigma$ , a power being repeated for states reached in the same number of time intervals from a fixed point in time.

These models can be elaborated in various ways, some of which are set out in [9, 12]. But I think I have said enough to indicate their usefulness in analyzing data on human stocks and flows and on the associated costs and benefits. They illustrate one way of analyzing stock-flow data but they do not provide the sole justification for the matrix framework of section 4 above. A framework for recording stocks and flows is needed since coherence and consistency are desirable features in statistics whatever the method of analysis.

## VI. LINKS BETWEEN SEQUENCES

As I have said, the proposal to divide up life into sequences is made with an eye to convenience in the regular reporting of statistics. As matters stand, stock statistics of student numbers are often collected by means of an annual questionnaire to schools; and it would not be difficult to collect at the same time flow statistics by adding a question about the position of each student a year ago, as is done in Holland [8]. It would be much more burdensome to require schools, as a matter of routine, to report personal and familial information about their students. This difficulty would not arise if the information were collected from the students and their families, as in the longitudinal study of Douglas and others in Britain [2, 3] and in the survey of Freytag and Wieszäcker in Baden-Württemberg [4, 5]. This source of information is an expensive one, however, and so, if it is used, anything like complete coverage cannot be expected. Of course, if statistics are collected by means of a linked system of compatible records or, better still, by a continuously updated, comprehensive system of individualized data, a discussion of sequence becomes largely irrelevant since the information in the vast, computerized data bank can be combined and extracted in any desired manner. But while these may be the methods of statistical collection of the future, they are not, with very limited exceptions, in operation at present, and so it makes sense to discuss the systematization of social statistics in terms of more familiar methods of collection.

The sequence approach is designed to describe different aspects of life in their own terms. For many purposes of organization and planning this may be quite adequate; for instance it enables us to trace the flow of students through the educational system and to work out the consequences of any expected changes in the structure of this system. If, however, we want to understand why different types of individual have very different experiences something more is needed. In terms of the sequence itself, we can only say that much depends on the kind of school attended and the qualifications obtained at different stages. If we want to probe deeper we must introduce into one sequence classifications characteristic of other sequences; for instance, we might introduce into the learning sequence such additional classifications as the social class of the student's

family, the urban or rural locality of that family and the early intelligence rating of the student himself. If we introduce such factors along with the type of school attended, we may hope to throw some light on such questions as whether at some stage in a student's career the influence of personal and familial characteristics wane and are replaced by the institutional characteristics of the school he attends.

Thus while for some purposes it is essential to combine information from different sequences, it is not clear that this requirement should be built into the regular reporting system. In the present example, the regular system can be regarded as a frame from which samples containing characteristics from many sequences could be drawn.

The models of the preceding section do not depend on the classifications used and so can be applied to mixed systems. However, as I have said, these models are not the only form of analysis and I have found that the method of regression on dummy variables is generally useful in dealing with social problems since many of the classifications in which we are interested are not numerical.

As an example suppose that we can classify the members of a particular vintage by sex, social class and intelligence rating and that we are interested in the effect of these influences on their scholastic performance. The different types of individual can be represented by the rows of a matrix consisting of 0's and 1's. For instance, with two sexes, two social classes and three ranges of intelligence, we should have  $2 \times 2 \times 3 = 12$  types. Adding a column of 1's to represent a general factor, the resulting matrix,  $X^*$  say, would take the form

$$(19) \quad X^* = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

In this matrix the twelve types are represented by the twelve rows. A 1 in the first column represents the general factor, in the second column being male, in the third column being female, in the fourth column being in the upper social class and so on.

With a complete matrix such as  $X^*$ , the product  $X^{*'}X^*$  is necessarily singular. This can be avoided if we delete one of the columns relating to each of the classifications. Let us agree to delete columns 2, 4 and 6 and denote the resulting matrix of type  $(12 \times 5)$  by  $X$ .

If the elements of a vector  $y$  denote the proportion of children in one of the twelve groups which succeeds in reaching a certain rung on the educational ladder or passing a certain educational test, then we can regress  $y$  on the columns of  $X$ . The regression equation is

$$(20) \quad y = Xb + e$$

where  $b$  denotes a vector of parameters and  $e$  denotes a vector of disturbances. The least-squares regression estimator,  $b^*$  say, of  $b$  is

$$(21) \quad b^* = (X'X)^{-1}X'y$$

The elements of  $b$  and  $b^*$  relate to columns 1, 3, 5, 7 and 8 of  $X^*$ . The coefficients relating to columns 2, 4 and 6 are all zero. Thus the calculated proportion of boys in the upper social class and the top intelligence range which passes the educational test in question is equal to the constant term and the regression coefficients measure additions to or subtractions from this number in respect of other characteristics. For example, the calculated proportion for girls in the lower social class and the lowest intelligence range is  $b_1 + b_3 + b_5 + b_8$  where  $b_1$  denotes the constant term,  $b_3$  denotes the differential for sex and so on. I shall give some examples of this type of analysis in the next section.

## VII. SOME NUMERICAL EXAMPLES

The following examples illustrate the analysis of social matrices by the methods described in the two preceding sections. The first three examples are based on matrices of order 114 which I shall not reproduce here since they are, in part, already available in [9, 12]. They are based on the scheme set out in my earlier paper [11] and not on the slightly different scheme embodied in the standard matrix of table 1 above. These matrices cover the first twenty years of life, that is ages 0 through 19, and in them the primary classification is by age, the secondary classification is by educational institution attended and the tertiary classification is by level of work. This last classification is effective only in the case of secondary schools where work is divided between that which precedes study for advanced-level examinations (not A-level) and study for these examinations (A-level).

The sixth example is based on an all-age (or age-free) matrix of order 44 relating the whole active sequence. The figures given below are based on the standard matrix: a condensed matrix of order 22 based on my earlier scheme is set out in [10]. In using matrices of different sizes there is an aggregation problem, the so-called problem of lumpability, which I have not yet investigated in this context. It is important, however, since, as Kemeny and Snell show in [7], an aggregated version of a system which can be interpreted as a Markov chain can only itself be so interpreted if certain conditions are satisfied.

Let us now turn to the examples.

(a) *The educational experience of two types of nineteen-year-old boy.* This example illustrates the backward model based on admission probabilities as set out in (4) of section 5(a) above. According to the assumptions of this model those in a given state at the end of a year come in fixed proportions from the

states that directly feed that state. Thus from the matrix inverse  $(I - G)^{-1}$  we can estimate the distribution over states in each of the first nineteen years of groups of individuals who in their twentieth year, that is at age nineteen, were in a given state. By still further reduction of the data we can form the following table.

TABLE 2  
TIME SPENT ON AVERAGE IN DIFFERENT ACTIVITIES IN THE  
FIRST TWENTY YEARS OF LIFE  
England and Wales, surviving male population at age 19, conditions  
of 1964-65

Activity	At University	Years Not in Full-time Education
Pre-school home	4.8	4.8
Nursery, primary and special schools	6.8	6.7
Secondary schools:		
Grammar	5.3	1.3
Others	1.3	3.1
Full-time further education n.e.s.	—	0.1
University	1.5	—
Left full-time education	0.4	4.0
Total	20.0	20.0

Note: Components do not always add up to totals because  
of rounding-off errors.

The figures in this table relate to average experience. Thus, since children do not normally change their secondary school (though some exception must be made when comprehensive schools are being formed), the table bears out the well-known fact that undergraduates tend to come from grammar schools and those who have left the educational system by age 19 tend to come from schools of a less academic character.

Instead of dividing secondary school experience by type of school we might as easily have divided it by level of work. Had we done so, we should have found that, out of 6.6 years at secondary school, undergraduates spent, on average, 4.7 years in not A-level work and 1.9 years in A-level work; whereas the corresponding figures for those outside the system of full-time formal education are 4.25 years and 0.15 years respectively.

All these figures are hypothetical in the sense that they reflect the average experience of two groups of boys who lived their lives under the conditions of 1964-1965. With changing admission probabilities, the experience of any particular vintage of the period would have been somewhat different.

(b) *The seven years from 13 through 19.* This example illustrates the forward model based on transition probabilities as set out in (6) of section 5(a) above. Suppose we ask the question: how much time will be spent in different activities after a certain age by individuals now in a given state? Table 3 answers this question for boys in eight initial states with respect to their average distribution

TABLE 3  
 TIME TO BE SPENT ON AVERAGE IN DIFFERENT EDUCATIONAL ACTIVITIES FROM AGE 13 THROUGH AGE 19  
 BY BOYS IN EIGHT SELECTED INITIAL STATES  
 England and Wales, conditions of 1964-65

		Years							
156	State	Age 0	Age 5: Nursery and Primary	Age 11: Primary	Age 13: Primary	Age 13: Secondary Modern	Age 13: Grammar	Age 13: Compre- hensive	Age 13: Other Normal Schools
	Activity								
	0. Not in full-time formal education	3.490	3,486	3,392	1,909	4,172	1,809	3,716	3,586
	1. Nursery and primary	0.009	0.009	0.029	1.032	0	0	0	0
	2. Secondary modern: (a) not A-level	1.349	1,366	1,277	0	2.552	0	0	0
	(b) A-level	0.006	0.006	0.006	0	0.011	0	0	0
	3. Grammar: (a) not A-level	0.789	0.799	0.924	2.235	0	3.309	0	0.006
	(b) A-level	0.324	0.328	0.384	1.260	0.019	1.302	0	0.003
	4. Comprehensive: (a) not A-level	0.444	0.449	0.508	0	0.044	0	2.869	0.323
	(b) A-level	0.048	0.049	0.051	0.004	0.014	0.004	0.238	0.089
	5. Other normal: (a) not A-level	0.193	0.195	0.115	0	0	0	0	2.607
	(b) A-level	0.010	0.010	0.006	0	0	0	0	0.139
	6. Special schools	0.057	0.017	0.009	0	0	0	0	0
	7. Further education n.e.s.	0.159	0.160	0.160	0.175	0.168	0.178	0.105	0.153
	8. Colleges of education	0.020	0.020	0.022	0.046	0.007	0.047	0.022	0.024
	9. Universities	0.014	0.106	0.119	0.336	0.016	0.347	0.052	0.066
	<b>Total</b>	<b>7.000</b>	<b>7.000</b>	<b>7.000</b>	<b>7.000</b>	<b>7.000</b>	<b>7.000</b>	<b>7.000</b>	<b>7.000</b>

Note: Components do not always add up to totals because of rounding-off errors.

of time during the seven years from the age of 13 through the age of 19, always on the assumption that the transition probabilities of 1964–1965 remain unchanged. This table is constructed from the entries in an inverse similar to  $(I - C)^{-1}$  but restricted to survivors, so that, apart from rounding-off errors, the entries in all the columns sum to 7.

The first column of Table 3 shows the average time that a batch of newborn boys could expect to spend in the various activities open to them between the ages of 13 and 19. Of this span of seven years, the first two are years of compulsory schooling, so that at most five out of the seven could be spent outside the system of full-time education. The entry in the first row and first column of the table is approximately 3.5, indicating that the average expectation at birth is for 1.5 years of voluntary full-time education between the end of compulsory school attendance and the age of 20.

If we compare the first two columns of the table we can see that the average expectations at birth are not much changed as we move on to 5-year-olds at normal schools. The prospects of these children are slightly better only because some of their mentally and physically handicapped contemporaries are at special schools.

A further improvement appears in the third column, relating to boys aged 11, the average age of transfer from primary to secondary school. At this age some children have already made the transfer, and those who have not yet done so have better prospects; the reason for this is that many of them are likely to be attending a higher type of primary school, the so-called preparatory school, which keeps children up to 13 and is specifically intended to prepare them for a grammar-school career. This point is made more evident in the fourth column: the boys who are still at primary school at 13 are almost all destined to enter grammar school.

The remaining four columns show the prospects of 13-year-olds at the four types of normal secondary school. The extreme contrast appears between the pupils of secondary modern schools and those of grammar schools: the former can expect, on average, less than a year of voluntary full-time education between the ages of 15 and 19, while the latter can expect more than three years of it over the same age-span. An equally marked difference can be seen in the average time spent by the two groups on A-level work and at institutions of higher education.

(c) *Boys and girls at secondary school.* Although the position is changing, girls have a smaller expectation than boys of reaching the highest rungs of the educational ladder, and it might be supposed that this was partly due to differential participation in secondary education. Table 4, which was constructed from the entries in matrices of the form  $(I - C)^{-1}$ , shows the secondary school experience of one thousand boys and one thousand girls under the conditions of 1964–1965.

This table shows that, on average, girls can expect to spend slightly longer at secondary school than boys can; and that, compared with boys, rather more time is spent at grammar schools and rather less at secondary modern schools. The difference between the sexes is that boys spend less time in work at the lower level but relatively much more time in advanced-level work, a moderate degree of success in which is required for admission to a university.

TABLE 4  
 EXPECTATION AT BIRTH OF TIME TO BE SPENT AT SECONDARY SCHOOL  
 BY ONE THOUSAND BOYS AND ONE THOUSAND GIRLS  
 England and Wales, conditions of 1964-65

Type of Secondary School	Boys			Girls		
	Not A-level	A-level	Total	Not A-level	A-level	Total
Secondary Modern	2,199	6	2,205	2,153	4	2,157
Grammar	1,116	315	1,431	1,267	232	1,499
Comprehensive	617	47	663	623	36	659
Other normal	335	10	345	338	6	344
Total	4,267	377	4,644	4,381	278	4,659

Note: Components do not always add up to totals because of rounding-off errors.

The table tells us about average times and throws no light on their distribution. How effective the greater time spent on advanced-level work by boys is can only be judged by looking at the qualifications achieved.

(d) *Determinants of educational performance.* This example relates to the method of regression on dummy variables described in the preceding section. The results shown are taken from a series of calculations made by Mrs. Mary Tuck on the basis of the data collected by Douglas and others as described in [2, 3]. I am indebted to Dr. Douglas of The M.R.C. Unit on Environmental Factors in Mental and Physical Illness and to the Medical Research Council for permission to use these data.

The Douglas survey traces the progression of a sample of about 5,000 children born in March 1946. In the early '60's these children were completing their lower level work at secondary school. With the data available they can be classified in a variety of ways which enables us to introduce personal and familial characteristics into the analysis of educational performance. In what follows, three characteristics are considered; sex, social class (two categories, *M* to denote middle class and *W* to denote manual working class) and ability and attainment rating at age eight, the earliest age available (three categories denoted by 1, 2, 3).

Let us now examine the probability of passing three successive educational hurdles:

- (i) passing O-level examinations, but without regard to number, subjects or marks;
- (ii) passing A-level examinations, but without regard to number, subjects or marks; and
- (iii) embarking on a first degree course, but without regard to subject.

In the first of these three cases a good fit is obtained with the form of equation given in section 6 above. In the second and third cases, however, this simple model, while fitting well enough for working-class boys and middle-class girls, does not fit well in the cases of middle-class boys and working-class girls.



According to the simple model, the probability that middle-class individuals in the first range of ability will pass at least one A-level examination or embark on a first degree course is greatly underestimated in the case of boys and greatly overestimated in the case of girls, while the reverse is true of individuals in the second and third ranges of ability. These defects can largely be removed by introducing two additional columns into the X-matrix, one of which has 1's in the second and third rows and the other of which has 1's in the eleventh and twelfth rows.

The regression estimates, their standard errors (in brackets) and a measure,  $\bar{R}^2$ , of goodness of fit are brought together in Table 5 below.

TABLE 5  
REGRESSION RESULTS IN ANALYSES OF EDUCATIONAL PERFORMANCE BY SEX,  
SOCIAL CLASS AND ABILITY RANGE

	O-level	A-level		First Degree	
	Simple Model	Simple Model	Extended Model	Simple Model	Extended Model
Constant term	0.813 (0.021)	0.431 (0.042)	0.528 (0.014)	0.244 (0.037)	0.323 (0.025)
Sex	-0.003 (0.019)	-0.051 (0.037)	-0.147 (0.013)	-0.065 (0.033)	-0.145 (0.023)
Social Class	-0.294 (0.019)	-0.166 (0.037)	-0.262 (0.013)	-0.095 (0.033)	-0.174 (0.023)
Ability range 2	-0.325 (0.023)	-0.217 (0.046)	-0.210 (0.013)	-0.127 (0.041)	-0.131 (0.023)
Ability range 3	-0.481 (0.023)	-0.259 (0.046)	-0.251 (0.013)	-0.126 (0.041)	-0.131 (0.023)
MB effect	—	—	-0.158 (0.018)	—	-0.111 (0.032)
WG effect	—	—	0.130 (0.018)	—	0.127 (0.032)
$\bar{R}^2$	0.98	0.82	0.99	0.62	0.91

In each column of Table 5 the constant term represents the calculated probability for middle-class boys in the first range of ability. The entries lower down in the column represent the amounts to be added or subtracted in respect of individuals in other groups. Thus the estimated probability that at least one O-level examination will be passed by working-class girls of the second range of ability is  $0.818 - 0.003 - 0.294 - 0.325 = 0.196$ ; and the corresponding estimated probability of embarking on a first degree course is, according to the extended model,  $0.323 - 0.145 - 0.174 - 0.131 + 0.127 = 0.001$  apart from a small rounding-off error.

The estimated probabilities are compared with the observed probabilities in Table 6 below.

A number of conclusions can be drawn from this analysis.

First, at the O-level stage there is no significant difference between the sexes but social class and ability both have large effects. The difference of effect between ability ranges 2 and 3 is large though not as large as the difference of effect between ability ranges 1 and 2.

TABLE 6  
OBSERVED AND ESTIMATED PROBABILITIES FOR VARIOUS TYPES OF INDIVIDUAL  
AT THREE STAGES IN THE EDUCATIONAL PROGRESSION

	O-level		A-level			First Degree		
	Observed	Simple Model	Observed	Simple Model	Extended Model	Observed	Simple Model	Extended Model
B M 1	0.833	0.818	0.540	0.432	0.528	0.352	0.244	0.323
B M 2	0.493	0.493	0.169	0.215	0.160	0.081	0.117	0.081
B M 3	0.282	0.337	0.109	0.173	0.118	0.082	0.118	0.082
B W 1	0.553	0.525	0.241	0.266	0.266	0.108	0.149	0.149
B W 2	0.186	0.200	0.062	0.049	0.056	0.035	0.022	0.018
B W 3	0.069	0.043	0.021	0.007	0.014	0.014	0.023	0.019
G M 1	0.804	0.815	0.382	0.381	0.381	0.163	0.179	0.179
G M 2	0.512	0.490	0.161	0.164	0.171	0.032	0.052	0.048
G M 3	0.364	0.334	0.127	0.123	0.031	0.052	0.052	0.048
G W 1	0.490	0.552	0.132	0.215	0.119	0.032	0.084	0.005
G W 2	0.189	0.196	0.034	-0.002	0.039	0	-0.043	0.001
G W 3	0.039	0.040	0.003	-0.043	-0.002	0.000	-0.042	0.002

Second, at the A-level and first degree stages, which are further removed from the ages of compulsory school attendance, a marked difference between the sexes emerges. At the same time the gap between ability ranges 2 and 3 gradually closes.

Third, at the two later stages the simple model does not account satisfactorily for the performance of middle-class boys or of working-class girls, but this defect is largely remedied by the extended model. For middle-class boys the difference of effect between those in ability range 1 and ranges 2 and 3 is substantially increased; while for working-class girls this difference is substantially diminished.

Finally, the models are based on the implicit assumption that the effects combine additively to determine the probabilities. The hypothesis that they combine multiplicatively was tested by repeating the regression calculations replacing the probabilities by their logarithms. The results were markedly less satisfactory.

(e) *Staying on at school.* A number of experiments have been made in projecting transition probabilities and similar measures such as the proportion of  $\lambda$ -year-olds, where  $\lambda$  exceeds the minimum school-leaving age, who are still at school. In this type of work a number of problems must be kept in mind.

In the first place, the necessary data are likely to be available for only a comparatively short period. The calculations described below are based on data for the fifteen years 1953 through 1967. Other relevant series are available for even shorter periods.

In the second place, while there is a tendency for more and more children to remain at school at any given age, the smooth development of this tendency is affected by various factors. Several instances can be given: in 1963, some

children who would otherwise have been able to leave school at the end of an autumn term were required to remain at school until the end of the following spring term; the growth of comprehensive schools is likely to encourage still more the existing tendency to remain at school longer; the school-leaving age was raised from 15 to 16 in 1972.

In the third place, the main purpose of making educational projections is to prepare for the expected situation some twenty or thirty years ahead since teacher-training programmes take time to work out and it is difficult to change teaching methods rapidly.

In the fourth place, and largely because of what has just been said, unbounded trends, such as linear or exponential trends, are dangerous. It seems desirable to work with sigmoid trends and these can be plausibly rationalized in terms of the simple epidemic model set out in (12) of section 5(e) above. However, in using this kind of trend allowance must be made for the sudden jumps already referred to.

In the fifth place, there are many plausible ways of fitting logistic trends. In the present context all of them give good fits as, indeed, do unbounded trends. But they give very different estimates of the two parameters,  $\beta$  and  $\gamma$ , and this is particularly important in the case of  $\gamma$ , the upper bound of the probability. As a consequence it is desirable to check the sensitivity of projections over the relevant range to the value of  $\gamma$  adopted and, if possible, to obtain an exogenous estimate of  $\gamma$ .

**TABLE 7**  
THE PROPORTION OF BOYS STILL AT SCHOOL AT AGES 16 THROUGH 19  
England and Wales

Age Group		1955	1960	1965	1970	1980	1990	$\infty$
16	Actual	0.183	0.228	0.292				
	Calculated (i)	0.181	0.234	0.293	0.359	0.663	0.732	0.794
	(ii)	—	—	—	—	—	—	—
	(iii)	0.181	0.232	0.293	0.362	0.677	0.777	1.000
	(iv)	0.181	0.233	0.293	0.360	0.667	0.750	0.850
17	Actual	0.094	0.130	0.162				
	Calculated (i)	0.096	0.127	0.165	0.210	0.317	0.428	0.682
	(ii)	0.095	0.130	0.165	0.194	0.231	0.246	0.253
	(iii)	0.096	0.127	0.165	0.208	0.305	0.395	0.558
	(iv)	0.096	0.127	0.165	0.210	0.318	0.432	0.700
18	Actual	0.037	0.051	0.060				
	Calculated (i)	0.038	0.051	0.063	0.070	0.078	0.080	0.081
	(ii)	0.038	0.051	0.062	0.071	0.080	0.083	0.085
	(iii)	0.039	0.050	0.062	0.074	0.090	0.098	0.105
	(iv)	0.039	0.050	0.063	0.078	0.114	0.151	0.250
19	Actual	0.005	0.005	0.008				
	Calculated (i)	0.005	0.006	0.007	0.008	0.010	0.011	0.011
	(ii)	0.005	0.006	0.007	0.009	0.011	0.012	0.015
	(iii)	0.005	0.006	0.007	0.009	0.011	0.012	0.015
	(iv)	0.005	0.006	0.007	0.009	0.012	0.015	0.030

These remarks are exemplified by the results given in Table 7 below. This table relates to the probabilities that boys of different ages from 16 through 19 will still be at school. The methods used in fitting the various regressions were devised by Mr. Allan Gordon, who also carried out the calculations. The results are tabulated here to three places of decimals only though three significant figures were used in fitting the equations.

If (12) is divided by  $c_{sr}$ , the relative change in the coefficient is expressed as a declining linear function of the level of the coefficient. Experience shows that this formulation leads to hopelessly inaccurate estimates of  $\beta$  and  $\gamma$ . Accordingly (12) was rewritten in continuous time and integrated to give

$$(22) \quad c_{sr} = \frac{\gamma}{1 + \exp(\alpha - \beta\gamma\theta)}$$

where  $\alpha \equiv \log_e[(\gamma/c_{sr}) - 1]$ . In this expression  $c_{sr}$  denotes the initial values of  $c_{sr}$  at time  $\theta = 0$ . Equation (22) can be written in the form

$$(23) \quad \log_e[(\gamma/c_{sr}) - 1] = \alpha - \beta\gamma\theta$$

and this equation forms the basis of the estimation methods employed.

In method (i) of Table 7, an initial value,  $\gamma_0$  say, of  $\gamma$  is assumed and (23) is used to estimate  $\alpha$  and  $\beta$ . The estimate of  $\alpha$  implies a new estimate,  $\gamma_1$  say, of  $\gamma$  and this value provides the basis for a second calculation. This iterative process converges and the final estimates of  $\beta$  and  $\gamma$  are adopted.

Method (ii) is similar to method (i) except that  $\gamma$  is obtained by ensuring that the median of the observed values of  $c_{sr}$  corresponds to a value of the dependent variable in (23) equal to the centroid of the regression line.

Method (iii) consists of working out the regression, (23) for many values of  $\gamma$  and choosing that value which maximizes  $r^2$ . If this value  $> 1$ , then the estimate  $\gamma = 1$  is adopted.

Method (iv) makes use of an exogenous estimate of  $\gamma$ .

In Table 7 the values of  $\gamma$  are shown in the final column. We can see from the table that we are dealing with a situation which is undergoing substantial changes so that the projected values of the coefficients are much greater than the values in the period of observation. In the case of the 16-year-olds, but not in the other cases, an allowance was made for a jump in 1972.

The calculated values of  $\gamma$  for the 16-year-olds are high and lie on either side of the exogeneous value. For the 17-year-olds, only the value of  $\gamma$  obtained by method (i) is at all similar to the exogenous estimate, the others being substantially lower. For the 18- and 19-year-olds the calculated values of  $\gamma$  are all substantially below the exogenous estimate.

This example illustrates the importance of trying to project coefficients as well as the difficulty of doing so. If anything like the projected changes comes about, a great deal of preparation will be needed. Moreover, these changes will affect not only the educational system but also labour relations and the organization of the labour market.

(f) *Life expectancies and their composition.* This example illustrates the use of the age-free matrix inverse  $(I - C)^{-1}$  for the calculation of life expectancies and their composition. The figures given below are taken from a table of the

kind given in [10] but recalculated in the form of the standard matrix used in this paper.

In this case the states of the initial stock-flow matrix of the active sequence are defined without regard to age. In it, individuals enter their pre-school home, stay in it for an unspecified time, move into the school system, progress through it until at some age they leave and are then classified by their leaving qualifications. After that they enter some form of further or higher education or take a job. Once in the labour force, most males remain in it, though they may return for a while to the educational system as students, until, if they survive long enough, they eventually retire.

A matrix of this kind was estimated for the male population of England and Wales in 1966. Since the population is not in stationary equilibrium, this matrix had to be adjusted; the methods used and their short-comings are discussed in [10]. From the adjusted matrix a *C*-matrix was calculated and from this a matrix inverse,  $(I - C)^{-1}$ , was derived.

If the estimates were accurate, this inverse would have the following properties.

- (i) The diagonal elements measure the mean time spent in a state by an individual about to enter that state.
- (ii) The off-diagonal elements measure the mean time spent in the state to which the row refers multiplied by the probability of reaching that state from the state to which the column refers.
- (iii) The sum of the elements in a column measures the expectation of life of an individual about to enter the state to which the column refers.
- (iv) From the life table we can discover the age at which a column sum is the expectation of life. This age is the average age at which individuals enter the state to which the column refers.

Although there are difficulties in combining demographic, educational and manpower statistics to give an accurate and consistent picture and although there are further difficulties in carrying out the adjustment to conditions of stationary equilibrium, the results obtained are reasonably reassuring. Thus, the male expectation of life at birth works out to 69.2 years compared with the official estimate of 68.5 years calculated by orthodox methods. The expectation of life on retirement is 10.4 years which implies that in 1966 British males retired on average at the age of 67. The years of retirement expected at birth are 6.64 and so  $6.64/10.42 = 0.64$  is the expectation at birth of surviving to retirement.

According to the matrix, the expectation of life at birth of 69.2 years is divided into the following partial expectancies; pre-school home, 5.1 years; education, 13.1 years; economic activity 44.4 years; and retirement 6.6 years. There is little doubt that the first two of these numbers are somewhat too high and that the third is somewhat too low. The picture is not very bad but more work is needed on the basic data and their adjustment.

As a final example of the use of this matrix let us see how the probability of going to a university changes with progression through the educational system. At birth the expectation is 0.08. If a boy leaves his secondary school with two or

more A-level certificates, the requirement for admission to a university, this probability rises to 0.60. If, on the other hand, he leaves without even any O-level certificates, the probability drops to 0.02. This figure may well be too high but, in the British system, it is at least possible to repair a disastrous school career at a later stage in an institution of further education; and a certain number of individuals avail themselves of this possibility.

(g) *Referrals in a psychiatric service system.* In [1], Baldwin gives a very interesting input-output table relating to the psychiatric service system of North-east Scotland which is centred on Aberdeen. This system is divided into nine branches as shown in Tables 8 and 9 below. The flows into, within and out of this system are termed "referrals"; a patient is referred into one of the branches of the system, or from one branch to another or out of the system altogether.

**TABLE 8**  
**DIRECT REFERRALS WITHIN THE SYSTEM PER 1,000 NEW ENTRANTS INTO EACH**  
**STATE OF THE PSYCHIATRIC SERVICE SYSTEM OF NORTH-EAST SCOTLAND IN 1965**  
**1,000 C**

From									
To	1	2	3	4	5	6	7	8	9
1. Out-patients	9	3	503	108	65	158	27	2	250
2. In-patients	205	57	248	585	130	357	459	173	638
3. Day-patients	32	26	17	17		7		15	28
4. Domiciliary visits					6				
5. Domiciliary treatment	18	35	17	17		11		54	28
6. Hospital consultations			2			2		1	
7. Other emergencies									
8. In-patient follow-up		305							28
9. Other psychiatric									
Total	264	426	785	729	201	535	486	246	972

The pattern of referrals within the system as set out in [1] is reproduced in Table 8. This can be interpreted as a  $C$ -matrix with each element multiplied by 1,000. Thus we can see from column 1 of Table 8 that, if 1,000 patients are referred into the system as out-patients, the direct effect is that 264 referrals will be made within the system and the balance, 736, will be referred out of the system. However, the 264 who are referred to another branch of the system, such as the 205 who are referred to in-patient treatment, will not all leave the system from that branch. As can be seen from column 2, over 30 per cent of those referred for in-patient treatment are next referred for in-patient follow-up. And of these, as can be seen from column 8, over 17 per cent are referred back for in-patient treatment. As a consequence, many indirect referrals will be made before the intake of 1,000 into any branch all succeed in getting out of the system. On the usual assumptions, the numbers can be calculated by forming the matrix inverse  $(I - C)^{-1}$ , as set out in Table 9.

Table 9 shows the direct and indirect consequences of 1,000 referrals from outside into any branch of the system. By combining Tables 8 and 9 we can see, for instance, that the entry of 1,033 in row 1 and column 1 of Table 9 can be decomposed into  $1,000 + 9 + 24 = 1,033$ . This means that the initial referral of 1,000 individuals to branch 1 of the system from outside generates 9 additional referrals directly and a further 24 referrals indirectly. From row 2 and column 1 of the two tables we can see that  $205 + 49 = 254$ , so that the initial referral of 1,000 individuals to branch 1 leads to 205 direct referrals and 49 indirect referrals to branch 2. Summing the entries in column 1 of Table 9 we obtain a figure of 1,440, indicating that if 1,000 patients are referred into the system at branch 1, 440 additional referrals will be made before the initial 1,000 have all left the system.

TABLE 9  
INITIAL, DIRECT AND INDIRECT REFERRALS PER 1,000 NEW ENTRANTS INTO  
EACH STATE OF THE PSYCHIATRIC SERVICE SYSTEM OF  
NORTH-EAST SCOTLAND IN 1965  
 $1,000(I - C)^{-1}$

From To	1	2	3	4	5	6	7	8	9
1. Out-patients	1,033	27	537	138	72	178	40	20	293
2. In-patients	254	1,149	423	711	170	456	534	216	819
3. Day-patients	42	37	1,048	44	8	27	18	23	64
4. Domiciliary visits	0	0	0	1,000	6	0	0	0	0
5. Domiciliary treatment	32	60	49	57	1,010	38	29	66	78
6. Hospital consultations	0	0	0	2	0	1,002	0	1	0
7. Other emergencies	0	0	0	0	0	0	1,000	0	0
8. In-patients follow-up	78	350	129	217	52	139	163	1,066	278
9. Other psychiatric	0	0	0	0	0	0	0	0	1,000
Total	1,440	1,624	2,186	2,169	1,318	1,841	1,784	1,391	2,533

Note: Components do not always add to totals because of rounding-off errors.

A similar interpretation can be put on the entries in the other columns of Table 9. But as in all such cases, these calculations are only strictly justified if the probabilities of movement from a state are independent of the path by which that state has been reached. It may well be that this condition is not satisfied with the states used in my example; and I must emphasize that it is only an example, even if the results are not wholly implausible. What is more important is that where a register has existed for some time it would be possible to test the assumption and to discover ways of defining states so that it is approximately satisfied. An example of how this might be done is given by Hall in [6].

### VIII. CONCLUSION

I have little to add by way of conclusion. I have tried to show that the division of life into sequences provides a manageable basis for the regular collection of

statistics on human stocks and flows; and, further, that the standard matrix provides a convenient method of arranging these statistics which is particularly useful from a purely technical point of view if they are collected by different agencies, as is frequently the case. Some of my examples are intended to show that there are many uses for information obtained from a single sequence, but there is no denying that other uses require the combination of classifications from two or more sequences. I see this as something it would be difficult to work into a regular reporting system but which could be dealt with from time to time using the regular system as a sampling frame.

Much of the paper is based on the assumption that countries are likely for the time being to make use of traditional methods of data collection over a large part of the field of social statistics. The division of life into sequences would, as I have said, be largely irrelevant with a continuously updated, comprehensive system of individualized data. I believe that in most countries we are a long way off such a system which would be expensive, technically difficult in practice and not without its political dangers. In the meantime, I believe that a great deal of useful information could be obtained from comparatively minor modifications of what I have called traditional methods.

#### IX. A LIST OF WORKS CITED

1. BALDWIN, J. A. *The Mental Hospital in the Psychiatric Service: A Case Register Study*. Oxford University Press, for the Nuffield Provincial Hospital Trust, 1971.
2. DOUGLAS, J. W. B. *The Home and the School*. MacGibbon and Kee, London, 1964.
3. DOUGLAS, J. W. B., J. M. ROSS and H. R. SIMPSON. *All Our Future*, Peter Davies, London, 1968.
4. FREYTAG, H. L., and C. C. VON WEIZSÄCKER. *Schulwahl und Schulsystem in Baden-Württemberg*. 2 vols., Heidelberg, 1968.
5. FREYTAG, H. L., and C. C. VON WEIZSÄCKER, (Eds.). *Schulwahl und Schulsystem: Modell theoretische Entwürfe-Verlaufs statistische Analysen*. 2 vols., Verlag Julius Beltz, Weinheim, 1969.
6. HALL, David J. Psychiatric prognosis in the middle-aged. In *Aspects of the Epidemiology of Mental Illness: Studies in Record Linkage*. Ed. J. A. Baldwin, Little Brown, Boston, in the press.
7. KEMENY, John G., and J. Laurie SNELL. *Finite Markov Chains*. D. Van Nostrand Co. Princeton, 1960.
8. NETHERLANDS CENTRAL BUREAU OF STATISTICS. *An Educational Matrix of the Netherlands for 1967*. N.C.B.S., 's-Gravenhage, 1969.
9. STONE, Richard. *Demographic Accounting and Model Building*. O.E.C.D., Paris, 1971.
10. STONE, Richard. "The fundamental matrix of the active sequence". In *Input-Output Techniques*. North-Holland, Amsterdam, 1972.
11. STONE, Richard and Giovanna, and Jane GUNTON. "An example of demographic accounting: the school ages". *Minerva*, Vol. VI, No. 2, 1968.
12. U.N. ECONOMIC AND SOCIAL COUNCIL. *An Integrated System of Demographic, Manpower and Social Statistics and its Links with the System of National Economic Accounts*. E/CN.3/394, mimeographed, 1970.