

TIME REVERSAL AND AGE DISTRIBUTIONS, I. DISCRETE-TIME MARKOV CHAINS

S. TAVARÉ,* *University of Sheffield*

Abstract

The connection between the age distribution of a discrete-time Markov chain and a certain time-reversed Markov chain is exhibited. A method for finding properties of age distributions follows simply from this approach. The results, which have application in several areas in applied probability, are illustrated by examples from population genetics.

MARKOV CHAINS; AGE DISTRIBUTION; DIFFUSION APPROXIMATION; GENETICS; TIME REVERSAL; DUALITY

1. Introduction

Several authors have considered statistical approaches to estimating the ‘age’ of a Markov process (Stigler (1970), Thompson (1976)). In this context, the ‘age’ is the absolute time origin of the process, and estimation is based on some ‘current’ observations. More recently, Levikson (1977) considered the problem of finding the age distribution of Markov processes with state space S ; a set of absorbing states, C , and transient states $B = S - C$. His method is to restart the process whenever C is visited by forcing instantaneous return to the nearest point of B . Between such returns, the probabilistic description of the two processes is identical. Given ‘current’ value $j \in B$, the age G_j has distribution determined by the limiting distribution of the time that has elapsed since C was last visited.

We now consider discrete-time Markov chains $\{X_n\}$, with corresponding return processes $\{\bar{X}_n\}$. We make a minor modification to Levikson’s restarting procedure. Instead of requiring instantaneous returns, we assume that it takes one step for a return to take place. This ensures that both processes have the same state space, which is a useful simplification (cf. Tavaré (1978a), Pakes (1978)). It is assumed that the return process is irreducible, aperiodic, and

Received 8 August 1978; revision received 13 December 1978.

* Present address: Department of Mathematics, University of Utah, Salt Lake City, UT 84112, U.S.A.

positive recurrent (as will usually be the case in genetic applications). The distribution of G_j is then given by

$$(1.1) \quad P(G_j = n) = g_j^{(n)} = \lim_{m \rightarrow \infty} P(\bar{X}_{m-n} \in C, X_{m-k} \notin C, 0 < k < n \mid \bar{X}_m = j).$$

Levikson used renewal theory methods to derive the distribution (1.1) in terms of the original chain, and some interesting properties of such an age. The purpose of this note is to give a simple method of deriving such results for an arbitrary restarting distribution. The idea is to relate the age distribution defined by (1.1) to a property of a certain time-reversed Markov chain. Such a representation considerably simplifies the analytical methods required to derive results about ages, and suggests some interesting problems in the theory of diffusion processes and weak convergence. Time-reversal has been used before in a genetic context to explain properties of asymptotic conditional distributions of the type

$$\lim_{m \rightarrow \infty} P(X_m = j \mid X_m \in B),$$

and

$$\lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} P(X_m = j, m < n \mid X_n \in B).$$

See, for example, Darroch and Seneta (1965), and Seneta (1966).

2. Two-barrier models

We suppose that $S = \{0, 1, \dots, M\}$, and $C = \{0, M\}$. This is the natural description of a wide variety of models that arise in population genetics. We begin with some notation. Suppose $\{X_n\}$ has transition matrix $P = (p_{ij})$, which we write in the form

$$(2.1) \quad P = \begin{bmatrix} 1 & \mathbf{0} & 0 \\ \mathbf{p}'_0 & Q & \mathbf{p}'_M \\ 0 & \mathbf{0} & 1 \end{bmatrix}.$$

Here, ' denotes transpose. The fundamental matrix is denoted $N = (n_{ij})$, where $n_{ij} = E$ (number of visits to $j \mid X_0 = i$); $i, j \in B$. Let $\pi_{ik} = P(X\text{-chain absorbed at } k \mid X_0 = i)$, $i \in B$, $k \in C$, and set $\pi_{00} = \pi_{MM} = 1$, $\pi_{0M} = \pi_{M0} = 0$. Let $\boldsymbol{\pi}_k = (\pi_{1k}, \pi_{2k}, \dots, \pi_{M-1,k})$.

Following Kemeny and Snell (1960), we have the results

$$(2.2) \quad \text{and} \quad \begin{aligned} N &= (I - Q)^{-1}, \\ \boldsymbol{\pi}'_k &= N\mathbf{p}'_k, \quad k \in C. \end{aligned}$$

To construct the transition matrix $\bar{P} = (\bar{p}_{ij})$ of the return process $\{\bar{X}_n\}$ we have to specify the elements $\{\bar{p}_{0l}, \bar{p}_{Ml} : l \in S\}$. The remaining elements of \bar{P} are identical to the corresponding elements of P , since probabilistic descriptions of the two processes are identical as long as the processes are in the set B . For notational convenience, we will set $\bar{p}_{0l} = r_{0l}, \bar{p}_{Ml} = r_{Ml}, l \in S$, and define

$$(2.3) \quad \text{and} \quad \begin{aligned} r_0 &= (r_{01}, r_{02}, \dots, r_{0, M-1}), \\ r_M &= (r_{M1}, r_{M2}, \dots, r_{M, M-1}). \end{aligned}$$

We can write \bar{P} in the form

$$(2.4) \quad \bar{P} = \begin{bmatrix} r_{00} & r_0 & r_{0M} \\ p'_0 & Q & p'_M \\ r_{M0} & r_M & r_{MM} \end{bmatrix}.$$

The n -step transition probabilities of X and \bar{X} are denoted $p_i^{(n)}$, and $\bar{p}_{ij}^{(n)}$, respectively, and the stationary distribution of the return process is denoted by $(\alpha_0, \alpha_1, \dots, \alpha_M)$.

It is then easy to check that for $0 < n \leq m$,

$$(2.5) \quad \begin{aligned} P(\bar{X}_{m-n} \in C, \bar{X}_{m-k} \notin C, 0 < k < n \mid \bar{X}_m = j) \\ = ({}_c\bar{p}_{0j}^{(n)} p_0^{(m-n)} / p_j^{(m)}) + ({}_c\bar{p}_{Mj}^{(n)} p_M^{(m-n)} / p_j^{(m)}), \end{aligned}$$

where $\{{}_c\bar{p}_{ij}^{(n)}\}$ are the n -step, C -avoiding transition probabilities of the return process, and $p_i^{(m)} = P(\bar{X}_m = i)$ is taken to be generated by any initial distribution for \bar{X}_0 . It follows from (1.1) and (2.5) that

$$(2.6) \quad g_j^{(n)} = {}_c\bar{p}_{0j}^{(n)} \frac{\alpha_0}{\alpha_j} + {}_c\bar{p}_{Mj}^{(n)} \frac{\alpha_M}{\alpha_j}.$$

One method of deriving the form of (2.6) is to relate the quantities on the right-hand side to properties of the original absorbing chain; for example, moments of the age distribution can be computed by matrix-generating function methods. However, it is clear that the limit derived in (2.6) is independent of the distribution of \bar{X}_0 . Hence we may assume that \bar{X}_0 has the stationary distribution, $(\alpha_0, \dots, \alpha_M)$; the chain $\{\bar{X}_n\}$ is then stationary. We may now reverse time, to obtain the dual chain $\{\bar{Y}_n\}$, which has the transition matrix $\bar{P}^* = (\bar{p}^*_{ij})$, where

$$(2.7) \quad \bar{p}^*_{ij} = \bar{p}_{ji} (\alpha_j / \alpha_i).$$

If we let $D_0 = \text{diag}\{\alpha_1, \dots, \alpha_{M-1}\}$, then it is easy to show that we can write

$$(2.8) \quad \bar{P}^* = \begin{bmatrix} r_{00} & \alpha_0^{-1} p_0 D_0 & \alpha_0^{-1} \alpha_M r_{M0} \\ \alpha_0 D_0^{-1} r'_0 & D_0^{-1} Q' D_0 & \alpha_M D_0^{-1} r'_M \\ \alpha_0 \alpha_M^{-1} r_{0M} & \alpha_M^{-1} p_M D_0 & r_{MM} \end{bmatrix}.$$

We also introduce the reversed absorbing process $\{Y_n\}$ which is derived from $\{\bar{Y}_n\}$ by making states in C absorbing. The corresponding transition matrix $P^* = (p^*_{ij})$ is

$$(2.9) \quad P^* = \begin{bmatrix} 1 & \mathbf{0} & 0 \\ \alpha_0 D_0^{-1} r'_0 & D_0^{-1} Q' D_0 & \alpha_M D_0^{-1} r'_M \\ 0 & \mathbf{0} & 1 \end{bmatrix}.$$

It should now be clear that the age, G_j , of the $\{X_n\}$ process is the absorption time of the reversed absorbing process, with $Y_0 = j$. By means of the representation of the age in terms of a well-specified Markov chain, we can now use standard Markov chain theory to derive a number of interesting properties of our age distribution. We begin with two lemmas.

Lemma 1. Let $\alpha = (\alpha_1, \dots, \alpha_{M-1})$. Then

$$\alpha_0 \propto r_M \cdot \pi'_0 + r_{M0},$$

$$\alpha_M \propto r_0 \cdot \pi'_M + r_{0M},$$

and

$$\alpha_j \propto (\alpha_0 r_0 N + \alpha_M r_M N)_j, \quad j \in B.$$

Proof. To find the stationary distribution of \bar{P} , we solve $(\alpha_0, \alpha, \alpha_M) = (\alpha_0, \alpha, \alpha_M) \bar{P}$. Using (2.4), write this as

$$(2.10) \quad \begin{aligned} \alpha_0 r_{00} + \alpha p'_0 + \alpha_M r_{M0} &= \alpha_0, \\ \alpha_0 r_{0M} + \alpha p'_M + \alpha_M r_{MM} &= \alpha_M, \end{aligned}$$

$$(2.11) \quad \alpha_0 r_0 + \alpha Q + \alpha_M r_M = \alpha.$$

From (2.11) and (2.2),

$$(2.12) \quad \alpha = \alpha_0 r_0 N + \alpha_M r_M N.$$

Substituting (2.12) into (2.10), and using (2.2) again yields

$$(2.13) \quad \alpha_0 (r_{0M} + r_0 \pi'_M) = \alpha_M (r_M \pi'_0 + r_{M0}).$$

The conclusion of the lemma follows by substituting (2.13) into (2.12).

In what follows, denote by $*$ the corresponding properties of the reversed absorbing chain $\{Y_n\}$. For instance, N^* is the fundamental matrix of the reversed absorbing process and T_j^* is the absorption time distribution of Y , given $Y_0 = j$.

Lemma 2. (i) $N^* = D_0^{-1}N'D_0$.

(ii) $\pi_k^{*'} = \alpha_k D_0^{-1}N'r'_k, \quad k \in C$.

Proof. Immediate from (2.2), (2.9).

To find the age distribution, we require the absorption-time distribution of the Y process. One can show that the n -step transition matrix of this process is given by

$$(2.14) \quad P^{**n} = \begin{bmatrix} 1 & \mathbf{0} & 0 \\ B_n^* p_0^{*'} & Q^{**n} & B_n^* p_M^{*'} \\ 0 & \mathbf{0} & 1 \end{bmatrix},$$

where $B_n^* = I + Q^* + \dots + Q^{*(n-1)}$. Then $P(T_j^* = n) = P(\text{age} = n \mid \text{'now' at } j) = g_j^{(n)}$. Let $\mathbf{g}_n = (g_1^{(n)}, \dots, g_{M-1}^{(n)})$. We then have the following result.

Lemma 3. For $n \geq 1$,

$$\begin{aligned} \mathbf{g}_n' &= Q^{*(n-1)}(p_0^{*'} + p_M^{*'}) \\ &= D_0^{-1}(Q')^{n-1}(\alpha_0 r_0' + \alpha_M r_M'). \end{aligned}$$

Proof. The first statement is just by definition of \mathbf{g}_n , and the final statement follows from (2.8).

The absorption probabilities of the reverse chain, Y , give the (limiting) probabilities of the X process having been restarted from a particular end of the state space, and will therefore be termed the restarting probabilities of X . A genetic interpretation, initially due to Levikson, will be given in the examples. We can now combine Lemmas 2 and 3 to give the following theorem.

Theorem 4. For $j \in B, n \geq 1$,

$$(i) \quad P(G_j = n) = g_j^{(n)} = \frac{\left(\sum_{k=0}^{M-1} r_{Mk} \pi_{k0}\right) \left(\sum_{k=1}^{M-1} r_{0k} p_{kj}^{(n-1)}\right) + \left(\sum_{k=1}^M r_{0k} \pi_{kM}\right) \left(\sum_{k=1}^{M-1} r_{Mk} p_{kj}^{(n-1)}\right)}{\left(\sum_{k=0}^{M-1} r_{Mk} \pi_{k0}\right) \left(\sum_{k=1}^{M-1} r_{0k} \pi_{kj}\right) + \left(\sum_{k=1}^M r_{0k} \pi_{kM}\right) \left(\sum_{k=1}^{M-1} r_{Mk} \pi_{kj}\right)}$$

(ii) $P(\text{last restart of } X \text{ process was at } M \mid \text{'now' at } j) = \pi_{jM}^*$

$$= \frac{\left(\sum_{k=1}^M r_{0k} \pi_{kM} \right) \left(\sum_{k=1}^{M-1} r_{Mk} n_{kj} \right)}{\left(\sum_{k=0}^{M-1} r_{Mk} \pi_{k0} \right) \left(\sum_{k=1}^{M-1} r_{0k} n_{kj} \right) + \left(\sum_{k=1}^M r_{0k} \pi_{kM} \right) \left(\sum_{k=1}^{M-1} r_{Mk} n_{kj} \right)}.$$

Proof. (i) follows from Lemma 3, and Lemma 1. (ii) follows from Lemma 2, and Lemma 1.

From the representation of the age of X as the absorption time of Y , we can readily deduce the moments of the age from standard Markov chain theory. Let $\mathbf{T} = (T_1, \dots, T_{M-1})$ be the vector of absorption times for the X process. Then, following Kemeny and Snell (1960), p. 49, we have

$$(2.15) \quad \begin{aligned} E\mathbf{T}' &= \mathbf{N}\mathbf{l}', \\ \text{Var } \mathbf{T}' &= (2\mathbf{N} - \mathbf{I})\mathbf{N}\mathbf{l}' - (\mathbf{N}\mathbf{l}')^2, \end{aligned}$$

where $\mathbf{l} = (1, 1, \dots, 1)$ and, for a vector $s = (s_1, s_2, \dots, s_{M-1})$, $s^2 = (s_1^2, s_2^2, \dots, s_{M-1}^2)$. The moments of the age distributions $\mathbf{G} = (G_1, \dots, G_{M-1})$ are just the moments of the absorption-time distributions $\mathbf{T}^* = (T_1^*, \dots, T_{M-1}^*)$, and these are computed using the following lemma.

Lemma 5.

$$E\mathbf{T}^{*\prime} = \mathbf{D}_0^{-1}\mathbf{N}'\mathbf{D}_0\mathbf{l}',$$

and

$$\text{Var } \mathbf{T}^{*\prime} = \mathbf{D}_0^{-1}(2\mathbf{N}' - \mathbf{I})\mathbf{N}'\mathbf{D}_0\mathbf{l}' - (\mathbf{D}_0^{-1}\mathbf{N}'\mathbf{D}_0\mathbf{l}')^2.$$

Another interesting problem related to the age process is the following. What is the (limiting) distribution of the number of visits to state $i \in B$, between successive visits to C , for a given 'current' position j ? This will give us added information on how the 'age process' is moving through its paths. We can easily reinterpret this distribution in terms of the reverse absorbing process. We are looking for the distribution of the number of visits to i before absorption, given $Y_0 = j$. Denote this random variable by N_{ji}^* . Then the distribution of N_{ji}^* is given by

$$(2.16) \quad \begin{aligned} P(N_{ji}^* = m) &= 1 - \frac{n_{ji}^*}{n_{ii}^*}, \quad m = 0 \\ &= \frac{n_{ji}^*}{n_{ii}^*} \cdot \frac{1}{n_{ii}^*} \cdot \left(1 - \frac{1}{n_{ii}^*}\right)^{m-1}, \quad m \geq 1. \end{aligned}$$

3. A class of restarting distributions of use in genetics

We consider the special case that arises from (2.3) by setting, for some $i \in B$,

$$(3.1) \quad \text{and} \quad \begin{aligned} r_{0j} &= 0, j \neq i; = 1, j = i, \\ r_{Mj} &= 0, j \neq M - i; = 1, j = M - i. \end{aligned}$$

The preceding results then specialise to:

$$(3.2) \quad \begin{aligned} \alpha_0 &\propto \pi_{M-i,0}, \\ \alpha_M &\propto \pi_{i,M}, \\ \alpha_j &\propto \pi_{M-i,0}n_{ij} + \pi_{i,M}n_{M-i,j}, \quad j \in T, \end{aligned}$$

$$(3.3) \quad g_j^{(n)} = \frac{\pi_{M-i,0}p_{ij}^{(n-1)} + \pi_{i,M}p_{M-i,j}^{(n-1)}}{\pi_{M-i,0}n_{ij} + \pi_{i,M}n_{M-i,j}}, \quad n \geq 1, \quad j \in B.$$

$$(3.4) \quad P(\text{last restart of } X \text{ process at } M \mid \text{'now' at } j) = \frac{\pi_{i,M}n_{M-i,j}}{\pi_{M-i,0}n_{ij} + \pi_{i,M}n_{M-i,j}},$$

and

$$(3.5) \quad n_{kj}^* = \frac{n_{jk}(\pi_{M-i,0}n_{ij} + \pi_{i,M}n_{M-i,j})}{(\pi_{M-i,0}n_{ik} + \pi_{i,M}n_{M-i,k})}.$$

The case $i = 1$ corresponds to the case considered by Levikson. The result of (3.2) suggests a new interpretation of a problem which has appeared in the genetics literature before. Diffusion processes are often used as approximations to the type of model described in the introduction. By suitably rescaling the time and state spaces, we arrive at a diffusion process $X(t)$ on $[0, 1]$, with drift and diffusion coefficients $m(x)$ and $v(x)$ respectively. In what follows, we will assume that

$$(3.6) \quad \begin{aligned} v(x) &= x(1-x), \\ m(x) &= x(1-x)\phi(x), \end{aligned}$$

where $\phi(x)$ is an arbitrary polynomial. The specification in (3.6) ensures that 0 and 1 are exit boundaries, corresponding to the absorbing barriers of the original model. Now write

$$(3.7) \quad G(x) = \exp \left\{ -2 \int^x \frac{m(y)}{v(y)} dy \right\}.$$

Then formal application of Wright's formula for stationary distributions (Wright (1937)) shows that X should have a stationary distribution $\alpha(x)$ of the form

$$(3.8) \quad \alpha(x) \propto \{G(x)v(x)\}^{-1}, \quad 0 < x < 1.$$

Since the specification (3.6) ensures that no stationary distribution can exist, what interpretation can be placed on (3.8)? This problem has received attention

before; for a fuller discussion, see Ewens (1963). We can now give another interpretation. The elements n_{ij} correspond to the pseudo-transient function of the diffusion, given by

$$(3.9) \quad \begin{aligned} n(p, x) &= \frac{2P_0(p)}{G(x)v(x)} \int_0^x G(y)dy, & 0 < x \leq p \\ &= \frac{2P_1(p)}{G(x)v(x)} \int_x^1 G(y)dy, & p \leq x < 1 \end{aligned}$$

where $P_1(p) = 1 - P_0(p) = P(X \text{ absorbed at } 1 | X(0) = p)$, and

$$(3.10) \quad P_1(p) = \int_0^p G(y)dy / \int_0^1 G(y)dy.$$

We introduce the instantaneous return process (IRP) $\bar{Y}(t)$, which is derived from X by jumping instantaneously to ε if the barrier at 0 is hit, and to $1 - \varepsilon$ if the barrier at 1 is hit. We may take $0 < \varepsilon < \frac{1}{2}$. This process is again a diffusion with state space $(0, 1)$. It clearly has a stationary distribution, $\alpha^*(x)$ and it can be shown from Feller (1954), p. 23 that it satisfies

$$(3.11) \quad \frac{1}{2} \frac{d}{dx} \{v(x)\alpha^*(x)\} - m(x)\alpha^*(x) = \begin{cases} D_1, & 0 < x < \varepsilon \\ D_2, & \varepsilon < x < 1 - \varepsilon \\ D_3, & 1 - \varepsilon < x < 1, \end{cases}$$

where D_i are constants, determined by requiring α^* to be continuous at ε and $1 - \varepsilon$, and integrable on $(0, 1)$. By analogy with (3.2), a substitution of the form

$$\alpha^*(x) \propto P_0(1 - \varepsilon)n(\varepsilon, x) + P_1(\varepsilon)n(1 - \varepsilon, x)$$

has the required property, and it follows that

$$(3.12) \quad \alpha^*(x) = \begin{cases} \frac{\left\{ \int_0^\varepsilon G(y)dy \right\}^{-1}}{G(x)v(x)} \int_0^x G(y)dy, & 0 < x \leq \varepsilon \\ \frac{1}{G(x)v(x)}, & \varepsilon < x < 1 - \varepsilon \\ \frac{\left\{ \int_{1-\varepsilon}^1 G(y)dy \right\}^{-1}}{G(x)v(x)} \int_x^1 G(y)dy, & 1 - \varepsilon \leq x < 1 \end{cases}$$

with

$$(3.13) \quad D_1 = \left\{ 2 \int_0^\varepsilon G(y)dy \right\}^{-1}, \quad D_2 = 0, \quad D_3 = \left\{ 2 \int_{1-\varepsilon}^1 G(y)dy \right\}^{-1}.$$

Wright's formula is derived by solving

$$(3.14) \quad \frac{1}{2} \frac{d}{dx} \{v(x)\alpha(x)\} - m(x)\alpha(x) = 0, \quad 0 < x < 1.$$

It follows by comparing (3.12), (3.13) with (3.14), that as $\varepsilon \rightarrow 0$, the stationary measure of the IRP is given by the solution of (3.14). A simple genetic explanation of such a process is as follows: we can suppose that rare mutation events reintroduce a particular allele into the population whenever it is lost. As $\varepsilon \rightarrow 0$, this corresponds to reintroducing a single copy of an allele at frequency $1/M$, or $1 - 1/M$, where M is the (large) population size.

4. Age distributions for a class of genetic models

We specialise the results of (3.1)–(3.5) to the case $i = 1$, in the context of a class of genetic models introduced by Cannings (1974). Consider a single locus in a population of fixed size M in each generation. There are two possible alleles, denoted A and B . Let X_n be the number of A -alleles at time n . In the absence of mutation pressures, $\{X_n\}$ is a Markov chain of the required type. Two models in this class are:

(i) Wright–Fisher model.

$$(4.1) \quad p_{ij} = \binom{M}{j} \left(\frac{i}{M}\right)^j \left(1 - \frac{i}{M}\right)^{M-j}, \quad i, j \in S.$$

(ii) Moran's model (Moran (1958)).

$$(4.2) \quad \begin{aligned} p_{i,i+1} &= p_{i,i-1} = \frac{i(M-i)}{M^2}, & i \in B, \\ p_{ii} &= 1 - p_{i,i+1} - p_{i,i-1}, & i \in B, \\ p_{ij} &= 0, & |i-j| > 1. \end{aligned}$$

These models both satisfy $\pi_{i,0} = 1 - i/M = 1 - \pi_{i,M}$. However, the elements of N are only of a manageable form in the case of Moran's model. (The recent paper of Piva and Holgate (1977) may produce an explicit, but complicated, expression for N in the case of the Wright–Fisher model.) For Moran's model one obtains

$$(4.3) \quad n_{ij} = \frac{M(M-i)}{(M-j)}, \quad 1 \leq j \leq i, \quad \text{and} \quad n_{ij} = M \frac{i}{j}, \quad i \leq j \leq M-1.$$

The elements of N for other members of the class are usually given by (rescaled) diffusion approximations (cf. Ewens (1969)). For the Wright model, we have

$$(4.4) \quad n_{ij} \approx \frac{2(M-i)}{(M-j)}, \quad 1 \leq j \leq i, \quad \text{and} \quad n_{ij} \approx 2 \frac{i}{j}, \quad i \leq j \leq M-1.$$

These are of course, just multiples of the exact result in (4.3). It follows from the previous formulae that, for the Moran model,

$$(4.5) \quad \pi_{jM}^* = \frac{n_{M-1,j}}{n_{M-1,j} + n_{1j}} = \frac{j}{M} = \pi_{jM},$$

$$ET_j^* = EG_j = \sum_{l=1}^{M-1} \frac{(n_{1l} + n_{M-1,l})n_{lj}}{n_{1j} + n_{M-1,j}} = \sum_{l=1}^{M-1} n_{jl} = ET_j.$$

These results could be anticipated by noticing that for the Moran model with restarting points 1 and $M - 1$, we have $\bar{P}^* \equiv \bar{P}$, i.e. the process is completely reversible, and so the age distribution of X is just the absorption distribution of X (cf. Watterson (1977)). The quantities π_{jM}^* have been interpreted by Levikson as the probability of the A allele being the oldest, given current frequency j . For the Moran model, $\pi_{jM}^* = \pi_{jM}$, as predicted by reversibility. However, for Wright's model, the process is not reversible, and so (4.5) is only the 'diffusion approximation' to what we want. It is therefore of interest to determine how accurate such an approximation is. In Table 1, exact results for the last restart probabilities, the mean and variance of the age, and mean and variance of the absorption times are compared with the corresponding diffusion approximations. It can be seen that the results are surprisingly good. For further comments on the adequacy of diffusion approximations to N in this context, see, for example, Tavaré (1979), Pollak and Arnold (1975).

5. One-barrier models

The state space of the X -chain is now $S = \{0, 1, \dots, M\}$, but $M = \infty$ is allowed. The absorbing state is $C = \{0\}$, and again $B = S - C$. The notation of the previous sections will be used again. In order for the return process to be recurrent, it is obvious that we must have the absorption probabilities of the absorbing process identically 1. This will be assumed throughout. We are again interested in the limiting distribution specified by (1.1). The transition matrix P of the X -chain is written

$$(5.1) \quad P = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{p}'_0 & Q \end{pmatrix}$$

while the return process transition matrix is given by

$$(5.2) \quad \bar{P} = \begin{pmatrix} r_{00} & \mathbf{r}_0 \\ \mathbf{p}'_0 & Q \end{pmatrix},$$

where $\mathbf{r}_0 = (r_{01}, r_{02}, \dots)$. \bar{P} is assumed to be the transition matrix of a positive-recurrent, irreducible, aperiodic Markov chain, whose stationary distribution we denote $(\alpha_0, \alpha_1, \dots)$. The limit in (1.1) is now given by

TABLE 1
Exact and diffusion approximations (DA) for Wright's model

<i>j</i>	<i>M</i> = 10			<i>M</i> = 15		
	π^*_{jM}	EG_j	Var G_j	π^*_{jM}	EG_j	Var G_j
	(DA)	(DA)	(DA)	(DA)	(DA)	(DA)
	ET_j	Var T_j		ET_j	Var T_j	
1	0.0771 (0.1000)	4.950 (5.658) 5.753	52.180 (62.880) 60.866	0.0505 (0.0667)	5.548 (6.503) 6.624	91.511 (114.313) 111.113
3	0.2918 (0.3000)	11.517 (11.279) 11.031	90.655 (90.119) 91.009	0.1921 (0.2000)	14.164 (14.070) 13.825	193.443 196.481 195.873
5	0.5000 (0.5000)	13.299 (12.913) 12.509	92.543 (90.487) 92.655	0.3297 (0.3333)	18.477 (18.134) 17.751	214.381 (212.541) 214.669
7	0.7082 (0.7000)	11.517 (11.279) 11.031	90.655 (90.119) 91.009	0.4659 (0.4667)	20.158 (19.761) 19.327	216.394 (212.966) 216.490
9	0.9229 (0.9000)	4.950 (5.658) 5.753	52.180 (62.880) 60.866	0.6022 (0.6000)	19.604 (19.225) 18.808	216.131 (213.215) 216.261
11				0.7388 (0.7333)	16.697 (16.443) 16.116	208.519 (208.541) 209.272
13				0.8701 (0.8667)	11.240 (10.863) 10.743	164.772 (169.467) 167.420

$$g_j^{(n)} = {}_0\bar{p}_{0j}^{(n)} \frac{\alpha_0}{\alpha_j}; \quad n \geq 1, \quad j \in B,$$

where ${}_0\bar{p}_{0j}^{(n)}$ is the n -step, 0-avoiding transition probability of the return process. Under the stated conditions, the limit is again independent of the distribution of \bar{X}_0 , which will now be taken to be $(\alpha_0, \alpha_1, \dots)$. The corresponding return processes $\{\bar{Y}_n\}$ and $\{Y_n\}$ have transition matrices given by

$$(5.3) \quad \bar{P}^* = \begin{pmatrix} r_{00} & \alpha_0^{-1} p_0 D_0 \\ \alpha_0 D_0^{-1} r'_0 & D_0^{-1} Q' D_0 \end{pmatrix},$$

and

$$(5.4) \quad P^* = \begin{pmatrix} 1 & \mathbf{0} \\ \alpha_0 D_0^{-1} \mathbf{r}'_0 & D_0^{-1} Q' D_0 \end{pmatrix}$$

respectively. Here $D_0 = \text{diag}\{\alpha_1, \alpha_2, \dots\}$, and $D_0^{-1} = \text{diag}\{1/\alpha_1, 1/\alpha_2, \dots\}$. The result concerning the fundamental matrix N is given by

$$(5.5) \quad N(I - Q) = (I - Q)N = I.$$

If $M < \infty$, then (5.5) reduces to (2.2). If $M = \infty$, then N is the minimal positive right (and left) inverse of $(I - Q)$. See, for example, Kemeny, Snell and Knapp (1966), p. 108. Again the age distribution of the X -process is related to the absorption time of the reverse process. The following lemma is useful in determining when the (recurrent) return process is positive, and in finding the stationary distribution in terms of N .

Lemma 6. (i) The return chain is positive if, and only if, $\sum_{k=1}^{\infty} r_{0k} E T_k < \infty$,

$$(ii) \quad \alpha_0 \propto 1, \quad \alpha_j \propto \sum_{k=1}^{\infty} r_{0k} n_{kj}, \quad j \in B.$$

Proof. Let $\alpha = (\alpha_1, \alpha_2, \dots)$. The chain is positive if, and only if, $(\alpha_0, \alpha) = (\alpha_0, \alpha) \bar{P}$ has a non-trivial non-negative solution, with $\alpha l' < \infty$, where $l = (1, 1, \dots)$. From (5.2), (5.5) we see that $\alpha = \alpha_0 r_0 N$ so that $\alpha l' < \infty$ if, and only if, $r_0 N l' < \infty$, i.e. $\sum_{k=1}^{\infty} r_{0k} E T_k < \infty$. (ii) follows immediately.

Theorem 7. (i) $N^* = D_0^{-1} N' D_0$.

$$(ii) \quad g_j^{(n)} = \frac{\sum_{k=1}^{\infty} r_{0k} p_{kj}^{(n-1)}}{\sum_{k=1}^{\infty} r_{0k} n_{kj}}, \quad n \geq 1, \quad j \in B.$$

$$(iii) \quad E T^{*'} = D_0^{-1} N' D_0 l'.$$

$$(iv) \quad \text{Var } T^{*'} = D_0^{-1} (2N' - I) N' D_0 l' - (D_0^{-1} N' D_0 l')^2.$$

Proof. These follow by arguments analogous to the equivalent results in the two-barrier case, with the help of Lemma 6.

Theorem 7 gives us the age distribution and some of its properties. If we specialise to the case $r_0 = (1, 0, \dots)$, so that the return is made to state 1, then we arrive at the equivalent of Levikson's models. In fact, we have the following corollary.

Corollary 8 (Levikson-type models). If $E T_1 < \infty$, then for $j, k \in B$,

$$(i) \quad g_j^{(n)} = P(G_j = n) = \frac{p_{1j}^{(n-1)}}{n_{1j}}, \quad n \geq 1.$$

$$(ii) \quad n_{jk}^* = \frac{n_{1k}n_{kj}}{n_{1j}}.$$

$$(iii) \quad EG_j = \sum_{k=1}^{\infty} \frac{n_{1k}n_{kj}}{n_{1j}}.$$

We can again use (2.15) to give the distribution of the number of visits to a state i for the reversed process for a given starting state j . The special case $i = j$ does not give us a good measure of the differences in behaviour of the two absorbing processes, since $n_{jj}^* = n_{jj}$. If $\bar{P} = \bar{P}^*$, then the two recurrent processes behave in an identical way. Then, of course, we will also have $N = N^*$. The role of such reversibility is discussed in the next section.

6. Reversibility and comments

The role of complete reversibility in the determination of properties of age distributions is apparent, since if $\bar{P} = \bar{P}^*$ the process is reversible, and then the two absorbing processes are identical. The idea of complete reversibility has been discussed in the case of certain diffusion results by Watterson (1977). In any case, we can derive age results in a simple way via time-reversal, and the limiting operation of (1.1) can then be given a simple intuitive meaning. In the one-barrier model, it is possible for the (recurrent) return process to be null, in which case the limit in (1.1) need not exist, and need not be independent of the initial distribution. Pakes (1978) discusses such chains in the case of Levikson-type return boundaries. However, if we are prepared to believe that the return process is stationary, which seems reasonable from the point of view of applications, then the results of Section 5 still follow. They will be identical to the results of Pakes in the cases where the strong ratio limit property holds. It is possible that time reversal may also give some results in the case of transient return chains.

For an application of reversibility to the infinite-allele Moran model, see Kelly (1977), and for another approach to age distributions see Sawyer (1977). The methods of this paper can be extended, under some restrictions, to continuous-time Markov chains. The results will be presented in Tavaré (1978b).

Acknowledgements

I should like to thank Geoff Watterson and Tony Pakes for helpful remarks on an earlier draft of this paper, and the referee, whose comments have considerably improved the presentation.

References

- CANNINGS, C. (1974) The latent roots of certain Markov chains arising in genetics. I. Haploid models. *Adv. Appl. Prob.* **6**, 260–290.
- DARROCH, J. N. AND SENETA, E. (1965) On quasi-stationary distributions in absorbing, discrete time, finite Markov chains. *J. Appl. Prob.* **2**, 88–100.
- EWENS, W. J. (1963) The diffusion equation and a pseudo-distribution in genetics. *J. R. Statist. Soc.* **B25**, 405–412.
- EWENS, W. J. (1969) *Population Genetics*. Methuen, London.
- FELLER, W. (1954) Diffusion processes in one dimension. *Trans. Amer. Math. Soc.* **77**, 1–31.
- KELLY, F. P. (1977) Some exact results for the Moran neutral-allele model. *Adv. Appl. Prob.* **9**, 197–201.
- KEMENY, J. G. AND SNELL, J. L. (1960) *Finite Markov Chains*. Van Nostrand, Princeton, NJ.
- KEMENY, J. G., SNELL, J. L. AND KNAPP, A. W. (1966) *Denumerable Markov Chains*. Van Nostrand, Princeton, NJ.
- LEVIKSON, B. (1977) The age distribution of Markov processes. *J. Appl. Prob.* **14**, 492–506.
- MORAN, P. A. P. (1958) Random processes in genetics. *Proc. Camb. Phil. Soc.* **54**, 60–71.
- PAKES, A. (1978) The age distribution of a Markov chain. *J. Appl. Prob.* **15**, 65–77.
- POLLAK, E. AND ARNOLD, B. C. (1975) On sojourn times at particular gene frequencies. *Genet. Res., Camb.* **25**, 89–94.
- PIVA, M. AND HOLGATE, P. (1977) The eigenvectors of a finite population model. *Ann. Hum. Genet.* **41**, 103–106.
- SAWYER, S. (1977) On the past history of an allele now known to have frequency p . *J. Appl. Prob.* **14**, 439–450.
- SENETA, E. (1966) Quasi-stationary distributions and time reversion in genetics. *J. R. Statist. Soc.* **B28**, 253–277.
- STIGLER, S. M. (1970) Estimating the age of a Galton–Watson branching process. *Biometrika* **57**, 505–512.
- TAVARÉ, S. (1978a) Age distributions for Markov processes in genetics (abstract). *Adv. Appl. Prob.* **10**, 17–19.
- TAVARÉ, S. (1978c) Time reversal and age distributions. II. Continuous-time Markov chains.
- TAVARÉ, S. (1979) Sojourn times for conditioned Markov chains. *Theoret. Popn Biol.* **15**, 108–113.
- THOMPSON, E. A. (1976) Estimation of age and rate increase of rare mutants. *Amer. J. Hum. Genet.* **28**, 442–452.
- WATTERSON, G. A. (1976) Reversibility and the age of an allele I. Moran's infinitely many neutral alleles model. *Theoret. Popn Biol.* **10**, 239–253.
- WATTERSON, G. A. (1977) Reversibility and the age of an allele II. Two allele models with selection and mutation. *Theoret. Popn Biol.* **12**, 179–196.
- WRIGHT, S. (1937) The distribution of gene frequencies in populations. *Proc. Nat. Acad. Sci. USA* **23**, 307–320.