

Ewens' sampling formula and related formulae: combinatorial proofs, extensions to variable population size and applications to ages of alleles

Robert C. Griffiths^{a,*}, Sabin Lessard^b

^a*Department of Statistics, University of Oxford, 1 South Parks Rd, Oxford OX1 3TG, UK*

^b*Département de Mathématiques et de Statistique, Université de Montréal, C.P. 6128, Succursale Centre-ville, Montréal, Qué., Canada H3C 3J7*

Received 30 September 2004

Available online 23 May 2005

Abstract

Ewens' sampling formula, the probability distribution of a configuration of alleles in a sample of genes under the infinitely-many-alleles model of mutation, is proved by a direct combinatorial argument. The distribution is extended to a model where the population size may vary back in time. The distribution of age-ordered frequencies in the population is also derived in the model, extending the GEM distribution of age-ordered frequencies in a model with a constant-sized population. The genealogy of a rare allele is studied using a combinatorial approach.

A connection is explored between the distribution of age-ordered frequencies and ladder indices and heights in a sequence of random variables. In a sample of n genes the connection is with ladder heights and indices in a sequence of draws from an urn containing balls labelled $1, 2, \dots, n$; and in the population the connection is with ladder heights and indices in a sequence of independent uniform random variables.

© 2005 Elsevier Inc. All rights reserved.

Keywords: Age distribution of alleles; Coalescent process; Ewens' sampling formula; GEM distribution; Infinitely-many-alleles model; Ladder indices and heights; Poisson–Dirichlet process; Urn model

1. Introduction

Ewens' (1972) sampling formula (ESF) is the probability distribution of the number of different types of genes and their frequencies at a selectively neutral locus under the infinitely-many-alleles model of mutation. The coalescent process of Kingman (1982) describing the genealogy of a sample underlies the sampling distribution. Kingman (1978) relates the sampling distribution to partition structures. The population model under which the ESF holds can be described as a diffusion process which contains a large class of finite-population

models in its limit domain of attraction as the population size tends to infinity and time is scaled appropriately. The class includes the Wright–Fisher model, the Moran model, and Cannings' (1974) exchangeable model generalizing the Wright–Fisher model. Earlier papers viewed the ESF as an approximate sampling formula in the Wright–Fisher model. The population gene frequencies are modelled as a diffusion process by Ethier and Kurtz (1981), and as a genealogical process by Griffiths (1980) and Donnelly and Tavaré (1987). Joyce and Tavaré (1987) relate the genealogical process to a linear pure birth and immigration process. Applied interest is that in the ESF the number of types is a sufficient statistic for the mutation rate. The stationary distribution of the population gene frequencies in the diffusion process model is the

*Corresponding author. Fax: +44 1865 272595.

E-mail addresses: griff@stats.ox.ac.uk (R.C. Griffiths), lessards@dms.umontreal.ca (S. Lessard).

Poisson–Dirichlet process. Kingman (1978) constructs the Poisson–Dirichlet process with a *paintbox* construction. Hoppe (1984) derived the ESF from an urn model representation of sampling genes. Donnelly (1986) studies a partition of sample genes by the age of their allele types in the urn model construction of Hoppe (1984). He shows an equivalence between the urn construction and Kingman’s *paintbox* construction. Donnelly and Kurtz (1996) study a particle process that relates the population frequencies, modelled as a measure-valued diffusion process, to the coalescent process. Combinatorial aspects of Ewens’ sampling formula, the Poisson–Dirichlet process, and age-ordered allele frequencies are detailed in Arratia et al. (2003). The distribution of non-mutant lines of descent after some given time in the past to the present time is studied by Griffiths (1980), Watterson (1984), Tavaré (1984) and Donnelly and Tavaré (1986) giving a generalization of the ESF to the distribution of allele types before and after the given time in the past. The distribution of age-ordered alleles in the ESF is derived in Donnelly and Tavaré (1986).

The age distribution of a mutation known to be of a given frequency in a population was first derived in a classic paper by Kimura and Ohta (1973). Recent papers studying the age distribution in a sample of genes or the population use genealogy in a coalescent process, a diffusion process, or a combination of the two (Slatkin and Rannala, 1997; Rannala and Slatkin, 1998; Wiuf and Donnelly, 1999; Stephens, 2000; Wiuf, 2000, 2001; Griffiths, 2003; Griffiths and Tavaré, 1998, 2003).

In this paper, we present an elementary proof of the ESF based on combinatorial arguments in the framework of the coalescent process. The approach is used to extend the sampling formula and the Poisson–Dirichlet distribution in the population to the case of a variable population size. The age-ordered distribution of gene frequencies in a sample and the population is also found in this case, extending known results.

The genealogy and age of a rare mutant type considered in Wiuf (2000, 2001) is studied using the combinatorial approach in this paper.

A connection is explored between the distribution of age-ordered frequencies and ladder indices and heights in a sequence of random variables. In a sample of n genes the connection is with ladder heights and indices in a sequence of draws from an urn containing balls labelled $1, 2, \dots, n$; and in the population the connection is with ladder heights and indices in a sequence of independent uniform random variables.

2. Ewens’ sampling formula: a combinatorial derivation

The ancestry of a random sample of n genes is described back in time by a coalescent tree, with vertices

where lineages have a common ancestor (Kingman, 1982). Mutations occur along the edges of the coalescent tree according to a Poisson process of intensity $\theta/2$. Many discrete population models are in the domain of attraction of the coalescent when time is measured appropriately. In a Wright–Fisher model of constant size N with an infinitely-many-alleles model of mutation, novel mutant types are formed at a rate of u per gene per generation. The distribution of the ancestral tree of a sample of n genes converges to the coalescent when time is measured in units of N generations, where $\theta = 2Nu$ is the scaled mutation rate per gene per generation and $N \rightarrow \infty$. The number of non-mutant ancestors of a sample of n genes is a death process back in time, where ancestral lines are lost by either mutation or coalescence. Griffiths (1980) and Tavaré (1984) study this death process in the entire population and in a sample of n genes. Ewens (1990) calls the events in the death process *defining events*. Lines lost by mutation determine the family tree of the mutation in the sample, with the number of members of the family as the number of leaves subtended by the mutation.

Label the sample genes and list them in the order in which they are lost backward in time, following either a mutation or a coalescence. In the case of coalescence, one of the two genes involved is chosen at random to continue back in time as a parent gene and the other gene is lost. There are $n!$ different ordered loss lists. If there are k different types of genes represented in the sample with n_l genes of type l for $l = 1, \dots, k$, then there are $n!/[n_1! \dots n_k!]$ possibilities for the positions of the genes of the different types in the list. However, if the types that have the same number of genes in the sample are not labelled, then this number of possibilities is divided by $[b_1! \dots b_n!]$, where b_j is the number of types represented j times in the sample for $j = 1, \dots, n$ such that $\sum_j b_j = k$.

Now consider the probability of each particular sequence of events. When i genes remain, the rate of mutation is $i\theta/2$ and the rate of coalescence is $i(i-1)/2$. The probability that a particular gene is the next one lost and that it is lost by mutation is $\theta/[i(\theta+i-1)]$. Similarly the probability that a particular gene of a given type is the next one lost and that it is lost by coalescence is $(j-1)/[i(\theta+i-1)]$, where j is the number of genes of the given type among the i remaining genes. Then it is clear that

$$n! \left(\frac{n!}{n_1! \dots n_k!} \right) \frac{1}{b_1! \dots b_n!} \frac{(n_1-1)! \cdot \theta \dots (n_k-1)! \cdot \theta}{1 \cdot \theta \dots n \cdot (\theta+n-1)} = \frac{n!}{1^{b_1} \dots n^{b_n}} \frac{1}{b_1! \dots b_n!} \frac{\theta^k}{\theta \dots (\theta+n-1)} \tag{1}$$

is the probability of having k types of genes with b_j types represented j times for $j = 1, \dots, n$ in a sample of size $n = \sum_j j b_j$. This is the sampling formula conjectured by

Ewens (1972) and proved by induction by Karlin and McGregor (1972).

Notice that, if the n genes are labelled, then the probability that n_l given genes among these are of type l for $l = 1, \dots, k$ is

$$\frac{n! \theta^k \prod_{l=1}^k (n_l - 1)!}{\prod_{i=1}^n i(\theta + i - 1)}.$$

Moreover, if the n sampled genes are labelled and the ancestry is traced back up to the point of m ancestral genes of given types, say types $1, \dots, m$, then it suffices to choose these m ancestral genes and proceed as above for the others to find

$$\frac{(n - m)! \theta^{k-m} \prod_{l=1}^m n_l! \prod_{l=m+1}^k (n_l - 1)!}{\prod_{i=m+1}^n i(\theta + i - 1)} \quad (2)$$

for the probability of having n_l given genes of type l for $l = 1, \dots, k$, types $1, \dots, m$ being ancestral and types $m + 1, \dots, k$ being mutant. This is the formula originally given by Watterson (1984), which also extends Kingman’s (1982) formula for the case $\theta = 0$.

3. Variable population size

If the population size is variable, the rate of mutation of i genes at time t back is $i\theta/2$ and the rate of coalescence is $i(i - 1)/[2\lambda(t)]$, where $\lambda(t) = N(t)/N(0)$ and t is expressed in units of $N(0)$ generations. Starting with n genes at the current time 0 and measuring time backward, the genes are lost by mutation or coalescence at random times $T_n < T_{n-1} < \dots < T_1$ according to a non-homogeneous death process of rate $i[\theta + (i - 1)\lambda(t)^{-1}]/2$, where i is the number of remaining genes at time t back. At time $T_i = t$, the probability that a particular gene is lost by mutation is $[\theta\lambda(t)]/[i(\theta\lambda(t) + i - 1)]$ and by coalescence $(j - 1)/[i(\theta\lambda(t) + i - 1)]$, where j is the number of remaining ancestral lineages of the particular gene type in the sample.

Assume n distinct genes at time 0 and allocate each of them a type such that n_l are of type l for $l = 1, \dots, k$. There are $n!/[n_1! \dots n_k!]$ possibilities. The genes are lost in order and there are $n!$ cases in all to consider. Decompose these cases by looking at the position of the last gene of each type defined as the number of remaining genes the last time there remains one gene of the given type. The outcome will be an ordered sequence $1 = i_{l_1} < \dots < i_{l_k} \leq n$, where i_{l_m} is the number of genes remaining just before the last gene of type l_m is lost, this type being the m th oldest, for $m = 1, \dots, k$, and (l_1, \dots, l_k) being a permutation of $(1, \dots, k)$. This sequence configuration is possible if and only if the inequalities $i_{l_m} \leq \sum_{v=1}^{m-1} n_{l_v} + 1$ hold for $m = 1, \dots, k$. The number of arrangements of the n genes satisfying these conditions is, starting from the last gene lost and ending

with the first one,

$$\prod_{m=1}^k n_{l_m} \frac{(\sum_{v=1}^m n_{l_v} - i_{l_m})!}{(\sum_{v=1}^m n_{l_v} - i_{l_{m+1}} + 1)!}, \quad (3)$$

with the convention that $i_{l_{k+1}} = n$, since there are n_{l_m} possible genes for the last one of type l_m to be lost and $(\sum_{v=1}^m n_{l_v} - i_{l_m})$ possible genes for the ones lost between the last one of type l_m and the last one of type l_{m+1} , for $m = 1, \dots, k$. The number of arrangements (3) can also be written as

$$\prod_{m=1}^k (n_{l_m}!) \binom{\sum_{v=1}^m n_{l_v} - i_{l_m}}{n_{l_m} - 1}.$$

The probability of each such sequence is

$$\mathbb{E} \left\{ \frac{\theta^k \prod_{l=1}^k [(n_l - 1)! \lambda(T_{i_l})]}{n! \prod_{i=1}^n [\theta \lambda(T_i) + i - 1]} \right\}.$$

Finally, if the types with the same number of genes are not labelled, the probability of having k types of genes with b_j types represented j times for $j = 1, \dots, n$ is

$$\frac{n! \theta^{k-1}}{(\prod_{l=1}^k n_l) (\prod_{j=1}^n b_j!)} \sum_{\mathbf{i}} a_{\mathbf{i}, \mathbf{l}} \mathbb{E} \left\{ \frac{\prod_{l=2}^k \lambda(T_{i_l})}{\prod_{i=2}^n [\theta \lambda(T_i) + i - 1]} \right\}, \quad (4)$$

where

$$a_{\mathbf{i}, \mathbf{l}} = \frac{\prod_{m=1}^k \binom{\sum_{v=1}^m n_{l_v} - i_{l_m}}{n_{l_m} - 1}}{\binom{n}{n_1, \dots, n_k}},$$

with $\mathbf{i} = (i_1, \dots, i_k)$ satisfying $1 = i_{l_1} < \dots < i_{l_k} \leq n$ and $i_{l_m} \leq \sum_{l=1}^{m-1} n_l + 1$, for $m = 1, \dots, k$, and $\mathbf{l} = (l_1, \dots, l_k)$ being a permutation of $(1, \dots, k)$. Note that

$$\sum_{\mathbf{i}} a_{\mathbf{i}, \mathbf{l}} = \prod_{m=1}^k \left\{ \frac{n_{l_m}}{\sum_{v=m}^k n_{l_v}} \right\}, \quad (5)$$

which is the probability that type l_m is the m th oldest, for $m = 1, \dots, k$. The product on the right side of (5) is obtained by conditioning on the older types in sequential order from $m = 1$ to $m = k$. The sum on the left side of (5) is obtained by partitioning the event according to the positions of the last genes of the k different types then using the above number of possible arrangements of all the genes for each case divided by $\prod_{m=1}^k (n_{l_m}!)$, in order not to distinguish genes within types. Also $\sum_{\mathbf{i}, \mathbf{l}} a_{\mathbf{i}, \mathbf{l}} = 1$. Notice that

$$\frac{n! \theta^{k-1}}{\prod_{l=1}^k n_l} \sum_{\mathbf{i}} a_{\mathbf{i}, \mathbf{l}} \mathbb{E} \left\{ \frac{\prod_{l=2}^k \lambda(T_{i_l})}{\prod_{i=2}^n [\theta \lambda(T_i) + i - 1]} \right\} \quad (6)$$

is the probability of having n_{l_m} genes of type l_m , this type being the m th oldest, for $m = 1, \dots, k$. In the case of a

constant population size, taking $l_k = k, \dots, l_1 = 1$, without loss of generality, this probability reduces to

$$\frac{(n-1)!}{n_k(n_k + n_{k-1}) \cdots (n_k + \cdots + n_2)} \frac{\theta^k}{\theta \cdots (\theta + n - 1)}, \tag{7}$$

which is the formula for the age-ordered types given by Donnelly and Tavaré (1986).

4. Ladder indices and heights in an urn model

There is a representation of the conditional distribution of the partial sums $\sum_{v=1}^m n_v$, $m = 1, \dots, k$, given \mathbf{i} as the distribution of ladder heights given ladder indices \mathbf{i} in an urn model.

The probability of a particular sequence \mathbf{i} is

$$\mathbb{E} \left\{ \frac{\prod_{c \notin \mathbf{i}} (c-1) \prod_{d \in \mathbf{i}} \theta \lambda(T_d)}{\prod_{i=1}^n [\theta \lambda(T_i) + i - 1]} \right\} \frac{(n-1)!}{\prod_{l=2}^k (i_l - 1)} \mathbb{E} \left\{ \frac{\prod_{i=1}^k \theta \lambda(T_{i_i})}{\prod_{i=1}^n [\theta \lambda(T_i) + i - 1]} \right\}. \tag{8}$$

The probability of an age-ordered configuration conditional on \mathbf{i} , from (6) and (8), is

$$h_i(\mathbf{n}) = \frac{n \prod_{l=2}^k (i_l - 1)}{\prod_{l=1}^k n_l} a_{i,1} = \frac{\prod_{l=2}^k (i_l - 1)}{(n-1)!} \prod_{m=1}^k \frac{(\sum_{v=1}^m n_v - i_m)!}{(\sum_{v=1}^{m-1} n_v - i_m + 1)!}. \tag{9}$$

By convention the factorial term in the denominator is taken as 1 when $m = 1$. Note that $\sum h_i(\mathbf{n}) = 1$, where summation is over $\sum_{l=1}^k n_l = n$ and $\sum_{l=1}^{m-1} n_l \geq i_m - 1$, $m = 2, \dots, k$.

Now label n balls in an urn $1, 2, \dots, n$ and draw out balls at random without replacement sequentially. Let $\mathbf{i} = (i_1, \dots, i_k)$ and $\mathbf{b} = (b_1, \dots, b_k)$ be ladder indices and ladder heights where record values occur in the numbers on the balls which are drawn. The last index k is defined such that $b_k = n$. The probability of a particular initial pair $i_1 = 1, b_1$ is n^{-1} . The distribution of ladder heights and indices in the urn model is equivalent to the distribution in a random permutation of $1, 2, \dots, n$.

After the $(m-1)$ th ladder index at i_{m-1} there are $n - i_{m-1}$ balls left in the urn, with $b_{m-1} - i_{m-1}$ balls having labels less than b_{m-1} . Thus for $m > 1$, the probability of obtaining i_m, b_m , conditional on (i_1, \dots, i_{m-1}) and (b_1, \dots, b_{m-1}) , is

$$\frac{b_{m-1} - i_{m-1}}{n - i_{m-1}} \cdots \frac{b_{m-1} - i_m + 2}{n - i_m + 2} \frac{1}{n - i_m + 1} = \frac{(b_{m-1} - i_{m-1})!}{(b_{m-1} - i_m + 1)!} \prod_{j=i_{m-1}}^{i_m-1} \frac{1}{n - j}. \tag{10}$$

The joint probability of a configuration \mathbf{i}, \mathbf{b} is therefore

$$P(\mathbf{i}, \mathbf{b}) = \frac{(n - i_k)!}{n!} \prod_{m=2}^k \frac{(b_{m-1} - i_{m-1})!}{(b_{m-1} - i_m + 1)!} = \frac{1}{n!} \prod_{m=1}^k \frac{(b_m - i_m)!}{(b_{m-1} - i_m + 1)!}. \tag{11}$$

The distribution of the age-ordered frequencies $\mathbf{n} = (n_1, \dots, n_k)$ is identical to the ladder height distribution with $b_m = \sum_{v=1}^m n_v$. This follows because $h_i(\mathbf{n})$ is proportional to $P(\mathbf{i}, \mathbf{b})$. Also by comparing the two distributions, we see that the marginal ladder index distribution is

$$P(\mathbf{i}) = \frac{1}{n \prod_{m=2}^k (i_m - 1)}. \tag{12}$$

5. Population frequencies

5.1. GEM distribution

In an infinite-population model where $\lambda(t) = 1$, $t \geq 0$, the age-ordered frequencies X_1, X_2, \dots have a GEM distribution given by

$$Z_1, Z_2(1 - Z_1), Z_3(1 - Z_2)(1 - Z_1), \dots, \tag{13}$$

where $\{Z_i, i \geq 1\}$ are mutually independent identically distributed (i.i.d.) random variables with density

$$\theta(1 - z)^{\theta-1}, \quad 0 < z < 1$$

(Donnelly and Tavaré, 1986; Ewens, 1990). This is a random partition representation (see Pitman (1996) and references therein). The unordered frequencies are distributed as a Poisson–Dirichlet point process (Kingman, 1978).

The population analogue of (2) is derived in Griffiths (1980), Donnelly and Tavaré (1986).

Let (n_1, \dots, n_k) be a sample taken from the GEM distribution arranged in age order and let $q(n_1, \dots, n_k)$ denote the age-ordered distribution defined by

$$q(n_1, \dots, n_k) = \sum_{1 \leq i_1 < i_2 < \dots < i_k} \frac{n!}{n_1! \cdots n_k!} \mathbb{E}(X_{i_1}^{n_1} \cdots X_{i_k}^{n_k}). \tag{14}$$

It is known that this distribution is given by (7), however, we give a short proof for completeness. Considering whether the oldest type in the sample is the oldest type in the population or not, we find that

$$q(n_1, \dots, n_k) = \binom{n}{n_1} \mathbb{E}(Z_1^{n_1} (1 - Z_1)^{n - n_1}) q(n_2, \dots, n_k)$$

$$\begin{aligned}
 &+ \mathbb{E}((1 - Z_1)^n)q(n_1, \dots, n_k) \\
 = &\binom{n}{n_1} \frac{\theta \Gamma(n_1 + 1) \Gamma(n - n_1 + \theta)}{\Gamma(n + \theta + 1)} q(n_2, \dots, n_k) \\
 &+ \frac{\theta}{n + \theta} q(n_1, \dots, n_k), \tag{15}
 \end{aligned}$$

where $q(n_2, \dots, n_k)$ is interpreted as 1 if $k = 1$. Simplifying (15) yields

$$\begin{aligned}
 q(n_1, \dots, n_k) &= \frac{\theta}{(n - n_1)} \frac{(n - 1)!}{(n - n_1 - 1)!} \\
 &\times \frac{\Gamma(n - n_1 + \theta)}{\Gamma(n + \theta)} q(n_2, \dots, n_k).
 \end{aligned}$$

and (7) follows by recurrence. The distribution of the age-ordered relative frequencies \mathbf{n}/n in a model with a constant-sized population, given by (7), converges to the GEM distribution as $n \rightarrow \infty$ because (7) is a sample distribution from the GEM distribution.

5.2. Variable population size

It is of interest to find the population distribution of age-ordered frequencies in a model with a variable-sized population. This extends the Poisson–Dirichet and GEM distributions. This distribution can be described in terms of: (i) the event times $\{T_j, j \geq 1\}$; (ii) a stochastic sequence \mathbf{i} generated by a mixture of Bernoulli trials $\{\chi_j, j \geq 1\}$ conditionally independent given $\{T_j, j \geq 1\}$ such that $P(\chi_j = 1 | T_j) = \theta \lambda(T_j) / [\theta \lambda(T_j) + j - 1]$, $P(\chi_j = 0 | T_j) = (j - 1) / [\theta \lambda(T_j) + j - 1]$; and (iii) the age-ordered frequencies, conditional on \mathbf{i} . The limit distribution for the age-ordered relative frequencies in a sample of genes conditional on \mathbf{i} is interpreted as the population distribution.

The sequence $\{T_j, j \geq 1\}$ is a reverse Markov chain, with transition distributions

$$\begin{aligned}
 P(T_j > t | T_{j+1} = s) &= \exp \left\{ -\frac{\theta}{2}(t - s)j - \binom{j}{2} \int_s^t \frac{du}{\lambda(u)} \right\} \\
 &= \exp \left\{ -\frac{j\theta}{2} \int_s^t \frac{du}{p_j(u)} \right\}, \tag{16}
 \end{aligned}$$

where $t > s$ and $p_j(u) = \theta \lambda(u) / [\theta \lambda(u) + j - 1]$. In the limit there is an entrance boundary at infinity in the process corresponding to an infinite population. It follows from (16) and

$$P(\chi_j = 1 | T_j = t) = p_j(t)$$

that

$$P(\chi_j = 1, T_j \in (t, t + dt) | T_j > t) = \frac{j\theta}{2} dt + o(dt).$$

The limit form of the distribution of the relative age-ordered sample frequencies $n_1/n, \dots, n_k/n$ as $n \rightarrow \infty$

with $n_i/n \rightarrow x_i$ for $i = 1, \dots, k$, is given by

$$\lim_{n \rightarrow \infty} n^{k-1} h_i(\mathbf{n}) = \prod_{l=1}^{k-1} \left(\sum_{j=1}^l x_j \right)^{i_{l+1} - i_l - 1}, \tag{17}$$

since

$$\begin{aligned}
 &\lim_{n \rightarrow \infty} n^{k-1} h_i(\mathbf{n}) \\
 &= \lim_{n \rightarrow \infty} \prod_{m=2}^k (i_m - 1) n^{i_m - 1} \frac{(\sum_{v=1}^m n_v - i_m)!}{(\sum_{v=1}^m n_v - 1)!} \\
 &\quad \times n^{-(i_m - 2)} \frac{(\sum_{v=1}^{m-1} n_v - 1)!}{(\sum_{v=1}^{m-1} n_v - i_m + 1)!} \\
 &= \prod_{m=2}^k (i_m - 1) \left(\sum_{v=1}^m x_v \right)^{-(i_m - 1)} \left(\sum_{v=1}^{m-1} x_v \right)^{i_m - 2} \\
 &= \prod_{l=2}^k (i_l - 1) \prod_{l=1}^{k-1} \left(\sum_{j=1}^l x_j \right)^{i_{l+1} - i_l - 1}. \tag{18}
 \end{aligned}$$

Note that $(n_1 - i_1)! = (n_1 - 1)!$ and $(\sum_{l=1}^k n_l - 1)! = (n - 1)!$ in simplifying the first line of (18).

The distribution mentioned in (iii) is related to the distribution of ladder heights and indices in a sequence of i.i.d. uniform random variables on $[0, 1]$, as shown in the next subsection.

5.3. Ladder indices and heights in a sequence of uniform random variables

Let $\{U_l, l \geq 1\}$ be a sequence of i.i.d. uniform random variables on $[0, 1]$ and $\{S_m, m \geq 1\}$ the ladder heights (record values) which occur in the sequence $\{U_l, l \geq 1\}$ at random ladder indices \mathbf{i} (with $i_1 = 1$) such that $S_m = U_{i_m}, m \geq 1$. Then the joint probability of the first k indices (i_1, \dots, i_k) and density of these ladder heights $\{S_m, 1 \leq m \leq k\}$ is, by direct argument,

$$\prod_{m=1}^{k-1} s_m^{i_{m+1} - i_m - 1}. \tag{19}$$

The marginal probability of obtaining the indices is

$$\int \prod_{m=1}^{k-1} s_m^{i_{m+1} - i_m - 1} \prod_{m=1}^k ds_m = \frac{1}{i_k} \prod_{m=2}^k \frac{1}{i_m - 1}, \tag{20}$$

where integration is over $0 < s_1 < s_2 < \dots < s_k < 1$. The conditional distribution of $\{S_m, 1 \leq m \leq k\}$ given the indices is thus

$$i_k \prod_{m=1}^{k-1} (i_{m+1} - 1) s_m^{i_{m+1} - i_m - 1}. \tag{21}$$

We rescale by making a transformation from (S_1, \dots, S_k) to

$$\eta_m = \frac{S_m}{S_k}, \quad 1 \leq m < k, \quad S_k.$$

The Jacobian determinant of the inverse transformation is S_k^{k-1} and the joint density is found to be

$$\begin{aligned} i_k \prod_{m=1}^{k-1} (i_{m+1} - 1) \eta_m^{i_{m+1} - i_m - 1} S_k^{\sum_{v=1}^{k-1} (i_{v+1} - i_v - 1)} S_k^{k-1} \\ = i_k \prod_{m=1}^{k-1} (i_{m+1} - 1) \eta_m^{i_{m+1} - i_m - 1} S_k^{i_k - 1}. \end{aligned} \tag{22}$$

Integrating with respect to $0 < s_k < 1$, the density of $\{\eta_m, 1 \leq m < k\}$ is

$$\prod_{m=1}^{k-1} (i_{m+1} - 1) \eta_m^{i_{m+1} - i_m - 1}, \tag{23}$$

which is identical to the density (17) of the partial sums $\{\sum_{v=1}^m X_v, 1 \leq m < k\}$. In the limit as $n \rightarrow \infty, k \rightarrow \infty, S_k \rightarrow 1$, and it follows that the distribution of $\{\sum_{v=1}^m X_v, m \geq 1\}$ given \mathbf{i} is identical to the distribution of ladder heights in $\{U_l, l \geq 1\}$ given that they occur at ladder indices \mathbf{i} .

Distribution (17) is simplified by making a transformation to independent exponential random variables

$$\mathbf{X} = (X_1, \dots, X_{k-1}) \rightarrow \mathbf{V} = (V_1, \dots, V_{k-1}),$$

where

$$\sum_{l=1}^m X_l = \exp\left\{-\sum_{l=m}^{k-1} V_l\right\}. \tag{24}$$

The Jacobian determinant of the inverse transformation is

$$(-1)^{k-1} \prod_{m=1}^{k-1} \exp\left\{-\sum_{l=m}^{k-1} v_l\right\} \tag{25}$$

and making the transformation in (17), the density of \mathbf{V} is

$$\prod_{m=1}^{k-1} (i_{m+1} - 1) \exp\{- (i_{m+1} - 1) v_m\}. \tag{26}$$

That is, V_1, \dots, V_{k-1} are i.i.d. exponential random variables with rates $i_2 - 1, \dots, i_k - 1$. In the limit as $k \rightarrow \infty, S_k \rightarrow 1$, and there is a representation

$$\begin{aligned} X_1 &= e^{-\sum_{l=1}^{\infty} V_l}, \\ X_m &= \sum_{l=1}^m X_l - \sum_{l=1}^{m-1} X_l \\ &= e^{-\sum_{l=m}^{\infty} V_l} (1 - e^{-V_{m-1}}), \quad m \geq 2. \end{aligned} \tag{27}$$

An equivalent representation to (27) is that, for $m \geq 1$,

$$-\log\left(\sum_{l=1}^m X_l\right) = \sum_{j=i_{m+1}}^{\infty} (j-1)^{-1} \chi_j W_j, \tag{28}$$

where $\{W_j, j > 1\}$ is a sequence of i.i.d. exponential random variables with rate parameter unity.

A third representation of (27) as a random partition is

$$X_m = \xi_{m-1} \prod_{l=m}^{\infty} (1 - \xi_l), \quad m \geq 1, \tag{29}$$

where $\{\xi_l, l \geq 0\}$ are mutually independent random variables, with $\xi_0 = 1$, and for $m \geq 2$, ξ_m having the density of

$$(i_{m+1} - 1)(1 - z)^{i_{m+1} - 2}, \quad 0 < z < 1.$$

Eq. (29) is obtained by setting, for $l > 1$, $\xi_l = 1 - e^{-V_l}$.

The distribution of the age-ordered allele frequencies X_1, X_2, X_3, \dots conditional on the number of non-mutant lineages being i_1, i_2, i_3, \dots when the allele types arose by mutation has now been shown to be equivalent to the distribution of $S_1, S_2 - S_1, S_3 - S_2, \dots$, where S_1, S_2, S_3, \dots are successive ladder heights (record values) in a sequence of i.i.d. uniform $[0, 1]$ random variables given to occur at ladder indices i_1, i_2, \dots . Representations (27)–(29) thus apply to age-ordered frequencies.

The mean values of the age-ordered frequencies, conditional on \mathbf{i} , from (27) are for $m \geq 1$

$$\mathbb{E}(X_m | \mathbf{i}) = \frac{1}{i_m} \prod_{l=m+1}^{\infty} \left(1 - \frac{1}{i_l}\right). \tag{30}$$

The unconditional mean frequencies can be partially found. We have

$$\begin{aligned} \mathbb{E}(X_1) &= \mathbb{E}[\mathbb{E}(X_1 | \mathbf{i})] \\ &= \mathbb{E}\left[\prod_{l=2}^{\infty} \left(1 - \frac{1}{i_l}\right)\right] \\ &= \mathbb{E}\left[\prod_{j=2}^{\infty} \left(1 - \frac{1}{j} \chi_j\right)\right] \\ &= \mathbb{E}\left[\prod_{j=2}^{\infty} \left(1 - \frac{\theta \lambda(T_j)}{j(\theta \lambda(T_j) + j - 1)}\right)\right]. \end{aligned} \tag{31}$$

A similar calculation shows that for $m \geq 1$

$$\mathbb{E}(X_m) = \mathbb{E}\left[\frac{1}{i_m} \prod_{j=i_{m+1}}^{\infty} \left(1 - \frac{\theta \lambda(T_j)}{j(\theta \lambda(T_j) + j - 1)}\right)\right]. \tag{32}$$

An alternative expression to (32) is

$$\mathbb{E}(X_m) = \mathbb{E}\left[\frac{1}{\theta \lambda(T_{i_{m+1}}) + i_m} \prod_{l=i_{m+1}}^{\infty} \frac{\theta \lambda(T_l) + l}{\theta \lambda(T_{l+1}) + l}\right]. \tag{33}$$

Eq. (33) is found by simplifying terms in the product of (32) and shifting the product index in the denominator by unity.

In the usual constant-sized population case when $\lambda(t) = 1, t > 0$, we have

$$\mathbb{E}(X_m) = \mathbb{E}\left[\frac{1}{\theta + i_m}\right]. \tag{34}$$

Remark. Convergence of the product in (33) needs justification. Let $\{\tau_l, l \geq 1\}$ be independent exponential random variables with rates $\{l(l + \theta - 1)/2, l \geq 1\}$. In the constant-sized population case (with notation T_l°), $T_l^\circ = \sum_{k=1}^l \tau_k$. As $l \rightarrow \infty$, $T_l^\circ, T_l \rightarrow 0$, and $T_l \sim T_l^\circ$ because $\lambda(0) = 1$. We assume here that $\lambda(t)$ is continuous at $t = 0$ and $|\lambda'(0)| < \infty$. As $l \rightarrow \infty$ the general term of the product satisfies

$$\begin{aligned} \frac{\theta\lambda(T_l) + l}{\theta\lambda(T_{l+1}) + l} &\approx \frac{\theta + l + 1 + T_l^\circ\lambda'(0)}{\theta + l + 1 + T_{l+1}^\circ\lambda'(0)} \\ &= \frac{\theta + l + 1 + (\tau_l + T_{l+1}^\circ)\lambda'(0)}{\theta + l + 1 + T_{l+1}^\circ\lambda'(0)} \\ &\approx 1 + \frac{\tau_l\lambda'(0)}{l} \\ &\approx 1 + \frac{2Y_l\lambda'(0)}{l^3}, \end{aligned} \tag{35}$$

where $\{Y_l, l \geq 1\}$ are i.i.d. exponential random variables with unit rate. The product converges because of the cubic term in the denominator in (35).

5.4. Laplace transforms

The Laplace transform of $-\log(X_1)$, conditional on $\{T_l, l > 1\}$, is

$$\begin{aligned} \mathbb{E}[e^{\phi \log(X_1)}] &= \mathbb{E}[X_1^\phi] \\ &= \mathbb{E}\left[\prod_{l=2}^{\infty} \frac{i_l - 1}{i_l - 1 + \phi}\right] \\ &= \mathbb{E}\left[\prod_{j=2}^{\infty} \left(1 - \frac{\phi\chi_j}{j - 1 + \phi\chi_j}\right)\right] \\ &= \prod_{j=2}^{\infty} \left(1 - \frac{\phi}{j - 1 + \phi} \frac{\theta\lambda(T_j)}{\theta\lambda(T_j) + j - 1}\right) \end{aligned} \tag{36}$$

$$= \prod_{j=2}^{\infty} [1 - \rho_j(\omega_j - 1)]^{-1}, \tag{37}$$

with notation $\beta_l = \theta\lambda(T_l) + l - 1, \rho_l = \theta\lambda(T_l)/(l - 1)$, and $\omega_l = (1 + \phi/\beta_l)^{-1}$ for $l > 1$.

The moments of X_1 can be found from (36), for $k = 0, 1, \dots$, by setting $\phi = k$ to obtain

$$\mathbb{E}(X_1^k) = \mathbb{E}\left[\prod_{j=2}^{\infty} \left(1 - \frac{k}{k + j - 1} \frac{\theta\lambda(T_j)}{\theta\lambda(T_j) + j - 1}\right)\right]. \tag{38}$$

A representation shown by (37) is

$$-\log(X_1) = \sum_{j=2}^{\infty} \gamma_j, \tag{39}$$

where $\{\gamma_j, j > 1\}$ are mutually independent random variables with Laplace transforms

$$\mathbb{E}(e^{-\phi\gamma_j}) = [1 - \rho_j(\omega_j - 1)]^{-1}, \quad j > 1.$$

The random variable γ_j has an atom at zero with probability $(1 + \rho_j)^{-1}$, and a continuous density

$$\begin{aligned} \sum_{l=1}^{\infty} \left(\frac{\rho_j}{1 + \rho_j}\right)^l \frac{1}{1 + \rho_j} \frac{\beta_j^l \gamma^{l-1}}{(l-1)!} e^{-\beta_j \gamma} \\ = \frac{\rho_j}{1 + \rho_j} \frac{\beta_j}{1 + \rho_j} e^{-\frac{\beta_j}{1 + \rho_j} \gamma} \\ = \frac{\rho_j}{1 + \rho_j} (j - 1) e^{-(j-1)\gamma}, \quad \gamma > 0. \end{aligned} \tag{40}$$

Of course $-\log(X_1)$ is continuous, which agrees with

$$P\left(\sum_{j=2}^{\infty} \gamma_j = 0\right) = \prod_{j=2}^{\infty} (1 + \rho_j)^{-1} = 0,$$

since the series diverges to zero, because ρ_j is asymptotic to $\theta/(j - 1)$. Note directly from (36) that $\gamma_j = \chi_j \kappa_j$, for $j > 1$, where $\{\chi_j, j > 1\}$ and $\{\kappa_j, j > 1\}$ are independent with $\{\kappa_j, j > 1\}$ mutually independent exponential random variables with rates $\{j - 1, j > 1\}$. The Laplace transform of $-\log(X_m)$, conditional on i_m and $\{T_l, l > 1\}$ is

$$\begin{aligned} \mathbb{E}[e^{\phi \log(X_m)}] \\ = \mathbb{E}\left[(1 - e^{-V_{m-1}})^{\phi} \prod_{j=m}^{\infty} e^{-\phi V_j}\right] \\ = \prod_{j=2}^{i_m} \left(1 + \frac{\phi}{j - 1}\right)^{-1} \prod_{j=i_m+1}^{\infty} [1 - \rho_j(\omega_j - 1)]^{-1}. \end{aligned} \tag{41}$$

The first product in (41) is obtained from

$$\begin{aligned} \mathbb{E}[(1 - e^{-V_{m-1}})^{\phi}] &= (i_m - 1) \int_0^{\infty} e^{-(i_m-1)v} (1 - e^{-v})^{\phi} dv \\ &= (i_m - 1) \int_0^1 y^{i_m-2} (1 - y)^{\phi} dy \\ &= (i_m - 1) \mathcal{B}(i_m - 1, \phi + 1) \\ &= \prod_{j=2}^{i_m} \frac{j - 1}{j - 1 + \phi}. \end{aligned}$$

The structure of (41) clearly implies that

$$-\log(X_m) = \sum_{k=2}^{i_m} \delta_k + \sum_{j=i_m+1}^{\infty} \gamma_j, \tag{42}$$

where $\{\delta_j, j > 1\}$ are independent exponential random variables such that δ_j has rate $j - 1$. Trying to simplify

(41) further by taking expectation with respect to i_m seems complicated.

As an application, it is of interest to calculate the probability p_O that the oldest type in a sample of genes is the oldest type in the population. In a constant-size population, from the GEM distribution

$$p_O = 1 - \mathbb{E}((1 - X_1)^n) = 1 - \frac{\theta}{\theta + n} = \frac{n}{\theta + n}.$$

In a model with a variable-sized population, using (38)

$$p_O = \sum_{k=1}^n (-1)^{k-1} \binom{n}{k} \times \mathbb{E} \left[\prod_{j=2}^{\infty} \left(1 - \frac{k}{k+j-1} \frac{\theta \lambda(T_j)}{\theta \lambda(T_j) + j - 1} \right) \right]. \quad (43)$$

5.5. GEM distribution equivalence with a constant population size

The consistency of representation (29) with the GEM distribution in a model with a constant-sized population where $\lambda(t) = 1, t \geq 0$, is shown in this section. A long calculation shows that moments in the finite-dimensional distributions of (29) coincide with moments in the GEM representation in the constant population size model. Let r_1, \dots, r_m be non-negative integers. In the GEM distribution

$$\begin{aligned} \mathbb{E} \left[\prod_{l=1}^m X_l^{r_l} \right] &= \prod_{l=1}^m \mathbb{E} \left[Z_l^{r_l} (1 - Z_l)^{\sum_{v=l+1}^m r_v} \right] \\ &= \theta^m \prod_{l=1}^m B \left(r_l + 1, \sum_{v=l+1}^m r_v + \theta \right) \\ &= \theta^m \prod_{l=1}^m \frac{r_l! \Gamma(\sum_{v=l+1}^m r_v + \theta)}{\Gamma(\sum_{v=l}^m r_v + \theta + 1)} \\ &= \frac{\Gamma(\theta)}{\Gamma(|r| + \theta)} \theta^m \prod_{l=1}^m \frac{r_l!}{\sum_{v=l}^m r_v + \theta}, \end{aligned} \quad (44)$$

where $|r| = \sum_{v=1}^m r_v$. To show the equivalence of the GEM distribution and (29), moments in (29) are calculated and shown to agree with (44). An identity that establishes the equivalence is that for $m > 1$

$$\begin{aligned} \mathbb{E} \left(\prod_{l=1}^m X_l^{r_l} \right) &= \frac{\theta}{\theta + r_m} \frac{r_{m-1}! r_m!}{(r_{m-1} + r_m)!} \\ &\times \mathbb{E} \left(\left[\prod_{l=1}^{m-2} X_l^{r_l} \right] X_{m-1}^{r_{m-1} + r_m} \right). \end{aligned} \quad (45)$$

It then follows by induction on m that (44) is satisfied. Note that for $m = 1$, directly from (29),

$$\mathbb{E}(X_1^{r_1}) = \mathbb{E} \left[\prod_{l=1}^{\infty} \frac{i_{l+1} - 1}{i_{l+1} + r_1 - 1} \right]$$

and for $m > 1$,

$$\begin{aligned} &\mathbb{E} \left[\prod_{l=1}^m X_l^{r_l} \right] \\ &= \mathbb{E} \left[\prod_{\alpha=1}^{m-1} \zeta_{\alpha}^{r_{\alpha+1}} (1 - \zeta_{\alpha})^{\sum_{v=1}^{\alpha} r_v} \prod_{l=m}^{\infty} (1 - \zeta_l)^{\sum_{v=1}^m r_v} \right] \\ &= \mathbb{E} \left[\prod_{\alpha=1}^{m-1} (i_{\alpha+1} - 1) B(r_{\alpha+1} + 1, i_{\alpha+1} + \sum_{v=1}^{\alpha} r_v - 1) \right. \\ &\quad \left. \times \prod_{l=m}^{\infty} \frac{i_{l+1} - 1}{i_{l+1} + |r| - 1} \right]. \end{aligned} \quad (46)$$

The expectation of the last product in (46) conditional on i_m is

$$\begin{aligned} &\mathbb{E} \left[\prod_{l=m}^{\infty} \frac{i_{l+1} - 1}{i_{l+1} + |r| - 1} \right] \\ &= \mathbb{E} \left[\prod_{j=i_m+1}^{\infty} \left(1 - \frac{\chi_j |r|}{j + |r| - 1} \right) \right] \\ &= \prod_{j=i_m+1}^{\infty} \left(1 - \frac{\theta |r|}{(j + |r| - 1)(j + \theta - 1)} \right) \\ &= \prod_{j=i_m+1}^{\infty} \frac{(j-1)(j + \theta + |r| - 1)}{(j + |r| - 1)(j + \theta - 1)} \\ &= \frac{1}{(i_m - 1)!} \frac{\Gamma(\theta + i_m) \Gamma(|r| + i_m)}{\Gamma(\theta + |r| + i_m)}. \end{aligned} \quad (47)$$

Simplification of the second last line in (47) to the last line follows by taking the limit of the product from $i_m + 1$ to n as $n \rightarrow \infty$ and applying Euler’s formula

$$\Gamma(z) = \lim_{n \rightarrow \infty} \frac{n! n^z}{\prod_{j=0}^n (z + j)}.$$

If $m = 1, i_m = 1$, this shows that

$$\mathbb{E}(X_1^{r_1}) = \frac{\Gamma(\theta + 1) \Gamma(r_1 + 1)}{\Gamma(\theta + r_1 + 1)},$$

in agreement with (44). The term containing i_m in (46), when $m > 1$ is thus

$$\begin{aligned} &(i_m - 1) B(r_m + 1, i_m + |r| - r_m - 1) \\ &\quad \times \frac{1}{(i_m - 1)!} \frac{\Gamma(\theta + i_m) \Gamma(|r| + i_m)}{\Gamma(\theta + |r| + i_m)} \\ &= \frac{\Gamma(r_m + 1) \Gamma(i_m + |r| - r_m - 1) \Gamma(\theta + i_m)}{(i_m - 2)! \Gamma(\theta + |r| + i_m)}. \end{aligned} \quad (48)$$

The probability that $i_m = i_{m-1} + j$, for $j \geq 1$, conditional on i_{m-1} is

$$\frac{\theta (i_{m-1} + j - 2)! \Gamma(\theta + i_{m-1})}{(i_{m-1} - 1)! \Gamma(i_{m-1} + j + \theta)}. \quad (49)$$

The expected value of the expression (48) conditional on i_{m-1} is obtained by multiplying (48) by (49) and

summing; that is

$$\frac{\theta\Gamma(r_m + 1)\Gamma(\theta + i_{m-1})}{(i_{m-1} - 1)!\Gamma(\theta + r_m + 1)} \times \sum_{j=1}^{\infty} B(\theta + r_m + 1, j - 1 + i_{m-1} + |r| - r_m) = \frac{\theta r_m! \Gamma(\theta + i_{m-1})}{(i_{m-1} - 1)!\Gamma(\theta + r_m + 1)} \times B(\theta + r_m, i_{m-1} + |r| - r_m). \tag{50}$$

Simplification from the second last line in (50) follows by expressing the Beta function as an integral, then summing in the integrand. The identity used is that for $a > 0, b > 0$,

$$\sum_{j=1}^{\infty} B(a + 1, b + j - 1) = B(a, b).$$

Multiplying the last line in (50) by the term containing i_{m-1} in (46)

$$(i_{m-1} - 1)B(r_{m-1} + 1, i_{m-1} + |r| - r_m - r_{m-1} - 1)$$

results in the expression

$$\frac{\theta}{\theta + r_m} \frac{r_m! r_{m-1}!}{(r_m + r_{m-1})!} \Gamma(r_{m-1} + r_m + 1) \times \frac{\Gamma(i_{m-1} + |r| - r_{m-1} - r_m - 1)}{(i_{m-1} - 2)!} \times \frac{\Gamma(\theta + i_{m-1})}{\Gamma(\theta + |r| + i_{m-1})}. \tag{51}$$

Comparing (48) and (51) establishes the identity (45) and therefore completes the proof of the equivalence of the GEM representation and (29).

6. Genealogy of a derived type in a population of constant size

In the case of a constant population size, the probability that a sample of n genes contains n_l genes of type l for $l = 1, \dots, k$ with $\sum_l n_l = n$ does not depend on the order in which the sampled genes are lost backward in time either by mutation or coalescence. Therefore, if a given type is represented r times and known to have been derived from another type in the sample, the probability for the last gene of this type to be lost when there remain $m + 1$ genes for $m = 1, \dots, n - r$ is given by

$$\frac{\binom{n - m - 1}{r - 1}}{\binom{n - 1}{r}}, \tag{52}$$

which converges to $q(1 - q)^{m-1}$ as n and r tend to infinity such that r/n converges to q . The time of

occurrence of this event, represented by T_{m+1} , is distributed as the sum of independent exponential variables of parameters $i(\theta + i - 1)/2$ for $i = m + 1, \dots, n$, whose expectation is

$$\sum_{i=m+1}^n \frac{2}{i(\theta + i - 1)} \tag{53}$$

Multiplying and summing over m , the mean age of the mutation that has given rise to the family of size r is

$$\sum_{m=1}^{n-r} \frac{\binom{n - m - 1}{r - 1}}{\binom{n - 1}{r}} \sum_{i=m+1}^n \frac{2}{i(\theta + i - 1)}. \tag{54}$$

The limit of (54) as $n \rightarrow \infty$ is

$$\frac{2q}{\theta - 1} \int_0^1 \frac{1 - v^{\theta-1}}{1 - v} \frac{v}{1 - (1 - q)v} dv. \tag{55}$$

If $\theta \rightarrow 0$ and $n \rightarrow \infty$, then the mean age, calculated directly from (54), is

$$\sum_{m=1}^{\infty} q(1 - q)^{m-1} \frac{2}{m} = -\frac{2q}{1 - q} \log(q). \tag{56}$$

Kimura and Ohta (1973) derived the classical formula (56). Griffiths and Marjoram (1996), Griffiths and Tavaré (1998), Wiuf and Donnelly (1999) and Stephens (2000) show that the mean age of a mutation that gave rise to a family of size r is (54) when $\theta = 0$. Griffiths (2003) shows that there is a simplification to

$$2r(n - r)^{-1} \sum_{j=r+1}^n j^{-1}. \tag{57}$$

In the treatment of the above authors only the lineages containing a given mutation are considered, with other mutations not affecting lineages.

Similarly to the derivation of (52), the probability for a gene of the derived type to be lost when there remain $m + 1$ genes among which j are of the derived type for $j = 2, \dots, r$ and $m = j, \dots, n - r + j - 1$ is

$$\frac{\binom{m - 1}{j - 1} \binom{n - m - 1}{r - j}}{\binom{n - 1}{r}},$$

whose limit is

$$\binom{m - 1}{j - 1} q^j (1 - q)^{m-j}$$

and the time of occurrence of the coalescence event responsible for this loss has expectation

$$\sum_{m=j}^{n-r+j-1} \frac{\binom{m-1}{j-1} \binom{n-m-1}{r-j}}{\binom{n-1}{r}} \times \sum_{i=m+1}^n \frac{2}{i(\theta+i-1)}, \tag{58}$$

whose limit as $n \rightarrow \infty$ is

$$\frac{2q^j}{\theta-1} \int_0^1 \frac{(1-v^{\theta-1})v^j}{(1-v)(1-(1-q)v)^j} dv. \tag{59}$$

Additionally, as $\theta \rightarrow 0$ the limit is

$$\sum_{m=j}^{\infty} \binom{m-1}{j-1} q^j (1-q)^{m-j} \frac{2}{m} = 2 \left(\frac{q}{1-q}\right)^j \int_q^1 \frac{(1-y)^{j-1}}{y^j} dy. \tag{60}$$

In the case $j = 2$, (60) evaluates to

$$\frac{2q}{1-q} + \frac{2q^2}{(1-q)^2} \log(q), \tag{61}$$

which corresponds to the expected time it takes for all genes of the derived type to coalesce.

The above treatment shows that, in the limit, the total number of genes remaining the first time there remain $j - 1$ genes of the derived type, denoted by $M(j)$, is distributed as a sum of j independent geometric variables of parameter q , and therefore the distribution of $qM(j)$ as q tends to 0 converges to the distribution of a sum of j independent exponential variables of parameter 1. Moreover, assuming $qM(j) = x$ fixed and multiplying the unit of time by q , the last time there remain j genes of the derived type converges in distribution to its mean, which is $2/x$, as θ and q tend to 0, since its variance, which is given by

$$\sum_{i=\frac{x}{q}+1}^n \frac{4}{q^2 i^2 (\theta+i-1)^2},$$

is bounded by

$$\int_{\frac{x}{q}-1}^{\infty} \left(\frac{4}{q^2 y^4}\right) dy = \frac{4q}{3(x-q)^3},$$

which converges to 0 as q tends to 0, in agreement with Wiuf (2000). This means that the last time there remain j genes of the derived type is distributed, in the limit, as twice the inverse of a gamma distribution.

7. Genealogy of a derived type with variable population size

When the population size is variable, the probability of having r genes of a derived type and $n - r$ genes of an ancestral type in a sample of size n for $r = 1, \dots, n - 1$ is

$$\frac{(n-1)! \theta^{n-r}}{r} \sum_{m=1}^{n-r} \frac{\binom{n-m-1}{r-1}}{\binom{n-1}{r}} \times E \left\{ \frac{\lambda(T_{m+1})}{\prod_{i=2}^n [\theta \lambda(T_i) + i - 1]} \right\}. \tag{62}$$

The probability of this event, given that there are two types in the sample, is proportional, as θ tends to 0, to

$$\frac{n}{r} \sum_{m=1}^{n-r} \frac{\binom{n-m-1}{r-1}}{\binom{n-1}{r}} E(\lambda(T_{m+1})), \tag{63}$$

which converges, as n and r tend to 0 such that r/n converges to q , to

$$L(q) = \sum_{m=1}^{\infty} (1-q)^{m-1} E(\lambda(T_{m+1})). \tag{64}$$

Moreover, given a frequency q of the derived type, the last gene of this type is lost by mutation when there remain $m + 1$ genes with probability

$$\frac{(1-q)^{m-1}}{L(q)} E(\lambda(T_{m+1})). \tag{65}$$

for $m \geq 1$, and the time of occurrence of this event is T_{m+1} with this probability.

The coalescent process in a variable-sized population can be coupled with a process in a population of constant size $N(0)$ by measuring time backwards in units of $\tau = \int_0^t \lambda(s)^{-1} ds$. In a population which decreases in size exponentially back in time $N(t) = N(0)e^{-\beta t}$, that is $\lambda(t) = e^{-\beta t}$, and $\beta t = \log(1 + \beta \tau)$. In such a case

$$E(\lambda(T_{m+1})) = E \left\{ \frac{1}{1 + \beta \mathcal{T}_{m+1}} \right\}, \tag{66}$$

where \mathcal{T}_{m+1} is distributed as a sum of independent exponential variables of parameters $i(i-1)/2$ for $i = m + 1, \dots, n$. Keeping $qm = x$ and $q\beta = b$ fixed as q tends to 0, the variable $\beta \mathcal{T}_{m+1}$ converges in distribution to $2b/x$. Then the variable $qM = X$, where M represents the number of genes remaining just after the loss of the last gene of the derived type has a limiting density function, as q tends to 0, that is proportional to

$$f(x) = \frac{e^{-x}}{1 + 2b/x} \tag{67}$$

for $x > 0$. Moreover, the time of occurrence of this event in time units of $qN(0)$ generations is distributed, as q tends to 0, as

$$(1/b) \log(1 + 2b/X). \quad (68)$$

Similarly, since the probability that the sampled genes are lost in a given order depends only on the position of the last gene of the derived type, the variable $qM(j) = X(j)$, where $M(j)$ represents the number of genes remaining the first time there remain $j - 1$ genes of the derived type, will be distributed, as q tends to 0, as a sum of j independent random variables, one of which has a density function proportional to $f(x)$ and the other $j - 1$ have an exponential distribution with parameter 1. Moreover, the time of occurrence of this event in time units of $qN(0)$ generations will be distributed, as q tends to 0, as

$$(1/b) \log(1 + 2b/X(j)). \quad (69)$$

Again, this is in agreement with (Wiuf, 2000, 2001).

References

- Arratia, A., Barbour, A.D., Tavaré, S., 2003. *Logarithmic Combinatorial Structures: A Probabilistic Approach*. European Mathematical Society Publishing House, Switzerland.
- Cannings, C., 1974. The latent roots of certain Markov chains arising in genetics: a new approach. I. Haploid models. *Adv. Appl. Probab.* 6, 260–290.
- Donnelly, P., 1986. Partition structures, Polya urns, the Ewens sampling formula, and the ages of alleles. *Theoret. Population Biol.* 30, 271–288.
- Donnelly, P., Kurtz, T.G., 1996. A countable representation of the Fleming-Viot measure-valued diffusion. *Ann. Probab.* 24, 698–742.
- Donnelly, P., Tavaré, S., 1986. The ages of alleles and a coalescent. *Adv. Appl. Probab.* 18, 1–19.
- Donnelly, P., Tavaré, S., 1987. The population genealogy of the infinitely-many neutral alleles model. *J. Math. Biol.* 25, 381–391.
- Ethier, S.N., Kurtz, T.G., 1981. The infinitely-many-neutral-alleles diffusion model. *Adv. Appl. Probab.* 13, 429–452.
- Ewens, W.J., 1972. The sampling theory of selectively neutral alleles. *Theoret. Population Biol.* 3, 87–112.
- Ewens, W.J., 1990. Population genetics theory—the past and the future. In: Lessard, S. (Ed.), *Mathematical and Statistical Developments of Evolutionary Theory*, NATO ASI Series C: Mathematical and Physical Sciences, vol. 299, Kluwer Academic Publishers, Dordrecht, The Netherlands, pp. 177–227.
- Griffiths, R.C., 1980. Lines of descent in the diffusion approximation of neutral Wright–Fisher models. *Theoret. Population Biol.* 17, 37–50.
- Griffiths, R.C., 2003. The frequency spectrum of a mutation, and its age, in a general diffusion model. *Theoret. Population Biol.* 64, 241–251.
- Griffiths, R.C., Marjoram, P., 1996. Ancestral inference from samples of DNA sequences with recombination. *J. Comput. Biol.* 3, 479–502.
- Griffiths, R.C., Tavaré, S., 1998. The age of a mutation in a general coalescent tree. *Stochast. Models* 14, 273–295.
- Griffiths, R.C., Tavaré, S., 2003. The genealogy of a neutral mutation. In: Green, P.J., Hjort, N.L., Richardson, S. (Eds.), *Highly Structured Stochastic Systems*. Oxford Statistical Science Series 27. Oxford University Press, Oxford, UK.
- Hoppe, F.M., 1984. Polya-like urns and the Ewens’ sampling formula. *J. Math. Biol.* 20, 91–94.
- Joyce, P., Tavaré, S., 1987. Cycles permutations and the structures of the Yule process with immigration. *Stochastic Process. Appl.* 25, 309–314.
- Karlin, S., McGregor, J.L., 1972. Addendum to a paper of W. Ewens. *Theoret. Population Biol.* 3, 113–116.
- Kimura, M., Ohta, T., 1973. The age of a neutral mutant persisting in a finite population. *Genetics* 75, 199–212.
- Kingman, J.F.C., 1978. Random partitions in population genetics. *Proc. R. Soc. London Ser. A* 361, 1–20.
- Kingman, J.F.C., 1982. The coalescent. *Stochastic Process. Appl.* 13, 235–248.
- Pitman, J., 1996. Random discrete distributions invariant under size-biased permutation. *Adv. Appl. Probab.* 28, 525–539.
- Rannala, B., Slatkin, M., 1998. Likelihood analysis of disequilibrium mapping, and related problems. *Amer. J. Hum. Genet.* 62, 459–473.
- Slatkin, M., Rannala, B., 1997. Estimating the age of alleles by use of intra-allelic variability. *Amer. J. Hum. Genet.* 60, 447–458.
- Stephens, M., 2000. Times on trees, and the age of an allele. *Theoret. Population Biol.* 57, 109–119.
- Tavaré, S., 1984. Line-of-descent and genealogical processes, and their application in population genetics models. *Theoret. Population Biol.* 26, 119–164.
- Watterson, G.A., 1984. Lines of descent and the coalescent. *Theoret. Population Biol.* 10, 239–253.
- Wiuf, C., 2000. On the genealogy of a sample of neutral rare alleles. *Theoret. Population Biol.* 58, 61–75.
- Wiuf, C., 2001. Rare alleles and selection. *Theoret. Population Biol.* 59, 287–296.
- Wiuf, C., Donnelly, P., 1999. Conditional genealogies and the age of a neutral mutant. *Theoret. Population Biol.* 56, 183–201.