



Within a sample from a population, the distribution of the number of descendants of a subsample's most recent common ancestor

John L. Spouge

Building 38A, Room 6N603, National Center for Biotechnology Information, Bethesda MD 20894, United States



ARTICLE INFO

Article history:

Received 26 June 2013

Available online 7 December 2013

Keywords:

Most recent common ancestor of a subsample
Coalescent theory

ABSTRACT

Sample n individuals uniformly at random from a population, and then sample m individuals uniformly at random from the sample. Consider the most recent common ancestor (MRCA) of the subsample of m individuals. Let the subsample MRCA have j descendants in the sample ($m \leq j \leq n$). Under a Moran or coalescent model (and therefore under many other models), the probability that $j = n$ is known. In this case, the subsample MRCA is an ancestor of every sampled individual, and the subsample and sample MRCAs are identical. The probability that $j = m$ is also known. In this case, the subsample MRCA is an ancestor of no sampled individual outside the subsample. This article derives the complete distribution of j , enabling inferences from the corresponding p -value. The text presents hypothetical statistical applications pertinent to taxonomy (the gene flow between Neanderthals and anatomically modern humans) and medicine (the association of genetic markers with disease).

© 2013 The Author. Published by Elsevier Inc.

This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/3.0/>).

1. Introduction

Consider the following hypothetical situation. Within a sample of n individuals, a subsample of m individuals share a morphological character. Upon genetic analysis, the m individuals share some genetic characters with a further $j - m \geq 0$ individuals within the sample. One might desire a p -value to test whether $j - m$ is “too small”, i.e., to test whether the concentration of the morphological character among individuals with the genetic characters is too excessive to reflect chance alone. This article derives a p -value by giving the sampling distribution of j . Depending on its context, a small p -value might suggest among other possibilities, e.g., that gene flow between the subpopulations represented by the subsample and its complement within the sample is not free (i.e., that the mathematical assumptions underlying the coalescent are violated), or that the genetic characters have a causal influence on the morphological character. The Discussion demonstrates how the p -value might be relevant to rejecting the hypothesis of free gene flow between Neanderthals and anatomically modern humans (Krings et al., 1997; Nordborg, 1998; Krings et al., 2000) or to associating a genetic disease or phenotype with a set of DNA markers necessary but not sufficient for it.

To determine the distribution corresponding to the p -value, consider Kingman's coalescent (Kingman, 1982a,b), where n individuals are sampled uniformly at random at time t_0 from a large population. Kingman examined a haploid population, but

coalescent models can also apply to sexual populations (Nordborg, 2004; Pollak, 2004; Wakeley et al., 2012). A pure death process D_t ($t \geq 0$) counts the ancestors of the sample at prior times $t_0 - t$. The process D_t transitions through the states $n \rightarrow n-1 \rightarrow \dots \rightarrow 2 \rightarrow 1$, with the state $D_t = k$ ($k = 2, \dots, n$) having a sojourn time τ_k exponentially distributed with parameter $d_k = \frac{1}{2}k(k-1)$, and with the state $D_t = 1$ absorbing.

The sample ancestry can be described using \mathcal{E}_n , the set of all equivalence relations on the n individuals. Consider the Markov chain $\mathcal{R}_n \rightarrow \mathcal{R}_{n-1} \rightarrow \dots \rightarrow \mathcal{R}_2 \rightarrow \mathcal{R}_1$, whose state-space is \mathcal{E}_n , where \mathcal{R}_k corresponds to having $D_t = k$ ancestors ($k = n, n-1, \dots, 1$). The variate \mathcal{R}_k partitions the n individuals into k equivalence classes, each equivalence class corresponding to an ancestor and containing the ancestor's descendants at time t_0 . Define the identity relation $\Delta = \{(i, i) : i = 1, 2, \dots, n\}$ and the trivial relation $\Theta = \{(i, j) : i, j = 1, 2, \dots, n\}$. Given $\xi, \eta \in \mathcal{E}_n$, let $\xi < \eta$ denote that η can be obtained from ξ by combining two equivalence classes in ξ , and in fact, $\Delta = \mathcal{R}_n < \mathcal{R}_{n-1} < \dots < \mathcal{R}_2 < \mathcal{R}_1 = \Theta$. The transition probabilities of the Markov chain $\{\mathcal{R}_k\}$ are

$$\mathbb{P}\{\mathcal{R}_{k-1} = \eta | \mathcal{R}_k = \xi\} = \begin{cases} 2/[k(k-1)] & \text{if } \xi < \eta \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Kingman shows that if ξ contains k equivalence classes,

$$\mathbb{P}\{\mathcal{R}_k = \xi\} = \frac{(n-k)!k!(k-1)!}{n!(n-1)!} \lambda_1! \lambda_2! \dots \lambda_k!, \quad (2)$$

E-mail addresses: spouge@ncbi.nlm.nih.gov, spouge@nih.gov.

where $\lambda_1, \lambda_2, \dots, \lambda_k$ are the sizes of the equivalence classes of ξ . Eq. (1) implies Eq. (2), so Eq. (2) holds for any model imposing Eq. (1) on the ancestry of a sample, in particular the Moran model (Moran, 1962; Kimura and Crow, 1964; Watterson, 1984; Donnelly and Tavaré, 1986a,b) (without mutation), or indeed any model of ancestry approximating a coalescent process closely enough.

Now, draw a subsample of m individuals uniformly at random from the sample of n individuals. The subsample has a most recent common ancestor (MRCA). For $1 \leq m \leq j \leq n$, let $p_{n,m;j}$ denote the probability that the subsample MRCA has j descendants within the sample. For $j = n$, e.g., the subsample has the same MRCA as the sample. From Theorem 2 in Saunders et al. (1984) with $l_1 = l_2 = 2$ (or Example 1 in Saunders et al. (1984)),

$$p_{n,m;n} = \frac{m-1}{m+1} \frac{n+1}{n-1}. \quad (3)$$

(See also p. 77 in Hein et al. (2005).)

In a standard notation (Graham et al., 1994), let $n^m = n(n-1) \cdots (n-m+1)$ denote the falling factorial for $1 \leq m \leq n$, with $n^0 = 1$. In addition to Eq. (3), we have the trivial boundary cases $p_{n,1;1} = p_{n,n;n} = 1$, so for $m < n$, consider the recursion

$$p_{n,m;m} = \frac{m(m-1)}{n(n-1)} p_{n-1,m-1;m-1} + \frac{(n-m)(n-m-1)}{n(n-1)} p_{n-1,m;m}, \quad (4)$$

which conditions on \mathcal{R}_2 , the two terms corresponding to coalescences: (1) within the subsample (probability $m(m-1)/[n(n-1)]$); and (2) outside of the subsample (probability $(n-m)(n-m-1)/[n(n-1)]$). Eq. (4) can provide an inductive proof of the formula

$$p_{n,m;m} = \frac{2(m-1)!}{(m+1)(n-1)^{m-1}} \quad (5)$$

from (Wiuf and Donnelly, 1999). (See also, e.g., p. 84 in Hein et al. (2005) and Eq. (1) in Rosenberg (2007).)

If $j = m$, Eq. (5) provides a p -value $p_{n,m;m}$, to test whether under the assumptions underlying the coalescent, subsample ancestries are likely to coalesce before coalescing with the remainder of the sample (see, e.g., p. 86 in Hein et al. (2005) for examples concerning Neanderthal ancestry (Nordborg, 1998 and Harris and Hey, 1999)). If $j > m$, then the relevant (left-sided) p -value becomes a sum $p_{n,m;j,\bullet} = \sum_{i=m}^j p_{n,m;i}$. With the motivating applications mentioned in the Introduction, Theorem 1 in the Results section extends the analytic formula for $p_{n,m;j}$ from $j = m$ and $j = n$ to $1 \leq m \leq j \leq n$.

2. Theory

Theorem 1. Let $1 \leq m \leq n$, and consider a sample whose ancestry satisfies Eq. (1). Under the set-up described above, for $m = 1$, definitions show that $p_{n,m;j}$ equals 1 if $j = 1$ and 0 otherwise. For $m > 1$,

$$p_{n,m;j} = \begin{cases} \frac{m-1}{m+1} \frac{2(j-2)^{m-2}}{(n-1)^{m-1}} & \text{for } 2 \leq m \leq j < n \\ \frac{m-1}{m+1} \frac{n+1}{n-1} & \text{for } 2 \leq m \leq j = n, \end{cases} \quad (6)$$

with $p_{n,m;j} = 0$ unless $2 \leq m \leq j \leq n$.

Remark. Eq. (6) reduces to Eq. (5) in the case $j = m$, as it should.

Proof. Note the following identity for $1 \leq a \leq b$:

$$a \sum_{i=a}^b (i-1)^{a-1} = \sum_{i=a}^b [i - (i-a)] (i-1)^{a-1} = \sum_{i=a}^b [i^a - (i-1)^a] = b^a \quad (7)$$

where the second equality follows because $i(i-1)^{a-1} = i^a$ and $(i-1)^{a-1}(i-a) = (i-1)^a$. Thus, $\sum_{j=m}^n p_{n,m;j} = 1$ for $2 \leq m \leq n$:

$$\begin{aligned} \frac{m-1}{m+1} \left[\frac{n+1}{n-1} + \sum_{j=m}^{n-1} \frac{2(j-2)^{m-2}}{(n-1)^{m-1}} \right] &= \frac{m-1}{m+1} \left[\frac{n+1}{n-1} + \frac{2}{(n-1)^{m-1}} \sum_{j=m}^{n-1} (j-2)^{m-2} \right] \\ &= \frac{m-1}{m+1} \left[\frac{n+1}{n-1} + \frac{2}{(n-1)^{m-1}} \frac{(n-2)^{m-1}}{m-1} \right] \\ &= \frac{m-1}{m+1} \left[\frac{n+1}{n-1} + \frac{2}{n-1} \frac{n-m}{m-1} \right] \\ &= 1, \end{aligned} \quad (8)$$

where the second equality follows from Eq. (7).

For $m = 1$, definitions yield $p_{n,1;1} = 1$, and for $m > 1$, $p_{n,m;j} = 0$ unless $m \leq j \leq n$. To set up an inductive proof of Theorem 1 for the cases in Eq. (6), let \mathcal{P}_i be the proposition that Theorem 1 holds for every $2 \leq m \leq j \leq n \leq i$. To start the induction, \mathcal{P}_2 is true, because by definition $p_{2,2;2} = 1$, agreeing with Eq. (6) for $2 \leq m \leq j = n \leq 2$ (the other case $2 \leq m \leq j < n \leq 2$ being vacuous).

For the inductive step, assume \mathcal{P}_{n-1} holds for some fixed $n \geq 3$. From Eq. (2) for $k = 2$, the probability that one of the two equivalence classes of \mathcal{R}_2 contains all m subsample individuals and has a total of i elements is

$$\left[\frac{(n-2)!2!1!}{n!(n-1)!} i!(n-i)! \right] \frac{(n-m)!}{(n-i)!(i-m)!} = \frac{2}{n-1} \frac{i^m}{n^m}, \quad (9)$$

because there are $(n-m)!/[(n-i)!(i-m)!]$ equally probable ways of forming the two equivalence classes of \mathcal{R}_2 by placing the m subsample individuals into an equivalence class of i elements.

As usual, let empty sums equal 0. To check that \mathcal{P}_n follows from \mathcal{P}_{n-1} , we check first that Eq. (6) holds for $2 \leq m \leq j < n$, then conclude from $\sum_{j=m}^n p_{n,m;j} = 1$ and Eq. (8) that Eq. (6) also holds for $2 \leq m \leq j = n$. For $2 \leq m \leq j < n$, then,

$$\begin{aligned} p_{n,m;j} &= \frac{2}{n-1} \sum_{i=j}^{n-1} \frac{i^m}{n^m} p_{i,m;j} \\ &= \frac{2}{n-1} \left[\frac{j^m}{n^m} \frac{m-1}{m+1} \frac{j+1}{j-1} + \sum_{i=j+1}^{n-1} \frac{i^m}{n^m} \frac{m-1}{m+1} \frac{2(j-2)^{m-2}}{(i-1)^{m-1}} \right] \\ &= \frac{m-1}{m+1} \frac{2}{n-1} \frac{1}{n^m} \left[j^m \frac{j+1}{j-1} + 2(j-2)^{m-2} \sum_{i=j+1}^{n-1} i \right] \\ &= \frac{m-1}{m+1} \frac{2}{n-1} \frac{1}{n^m} \\ &\quad \times \left\{ j^m \frac{j+1}{j-1} + 2(j-2)^{m-2} \frac{1}{2} [n^2 - (j+1)^2] \right\} \\ &= \frac{m-1}{m+1} \frac{2(j-2)^{m-2}}{(n-1)^{m-1}}, \end{aligned} \quad (10)$$

where the first equality is justified as follows. One of the two equivalence classes of \mathcal{R}_2 must contain all m individuals from

Download English Version:

<https://daneshyari.com/en/article/6372377>

Download Persian Version:

<https://daneshyari.com/article/6372377>

[Daneshyari.com](https://daneshyari.com)