## Precision of the estimation of S(p²)

If $a_i$ is any frequency observed in a multinomial distribution and $p_i$ the probability of any observation falling in that class, then it is well known and easily demonstrated that the expectations of the following simple functions of the observations are as shown below.

Function            Expectation

$$\Sigma(a) \qquad n\,\Sigma(p) \qquad = \qquad n$$
$$\Sigma\{a(a-1)\} \qquad n(n-1)\Sigma(p^2) \qquad = \qquad n(n-1)s_2$$
$$\Sigma\{a(a-1)(a-2)\}n(n-1)(n-2)\Sigma(p^3) \qquad = \qquad n(n-1)(n-2)s_3$$

$$\left.\right\} (1)$$

and so on.

In the absence of knowledge of the number, or true frequencies of the distinguishable classes, an unbiased estimate can be made of the parameter

$$\propto = \Sigma(p^2) = s_2 \tag{2}$$

by using the observable statistic

$$\Sigma\{a(a-1)\} \div n(n-1) \; ; \tag{3}$$

it may be required to assign a sampling variance of this estimate.

Now, the mean value of the square

$$\Sigma^2\{a(a-1)\} \tag{4}$$

may be expressed as that of the sum for all classes of

1 - 5 + 6
4 - 8
+ 2

$$a^2(a-1)^2$$
$$= a(a-1)(a-2)(a-3) + 4a(a-1)(a-2) + 2a(a-1) \tag{5}$$

of which the mean is

$$n(n-1)\left\{ 2s_2 + 4(n-2)s_3 + (n-2)(n-3)s_4 \right\} , \tag{6}$$

and that of the sum of products for all pairs of classes

$$2a(a-1)\,b(b-1), \tag{7}$$

of which the mean is

$$n(n-1)(n-2)(n-3)\,\Sigma\Sigma(2p^2q^2) , \tag{8}$$

and which together with the last term in $s_4$ gives

$$n(n-1)(n-2)(n-3)s_2^2 \tag{9}$$

Consequently, the sampling variance of the estimate of the parameter $\propto$ is

$$\left\{ 2s_2 + 4(n-2)s_3 - (4n-6)s_2^2 \right\} / n(n-1) \tag{10}$$

- X -

which is not expressible wholly in terms of $s_2$, but requires also a knowledge of the sum of the third powers, $s_3$. In general, we may expect $s_3$ to lie between $\alpha$ and $\alpha^2$, but the limits of uncertainty are generally rather wide, $\left\{ \frac{2\alpha(1-\alpha_2)}{2\alpha(1-\alpha)} \frac{\alpha_3 - (4\alpha_2-1)\alpha_2(1-\alpha_2)}{(\alpha_2(4\alpha_2-6)\alpha(1-\alpha)} \right\}$.

In the case that the frequencies are those of individual alleles in a self-sterility series, the probabilities p have a distribution, which, for statistical equilibrium, is known in terms of the parameter $\alpha$, and the population size, N.

The distribution is

$$\frac{1}{p} \cdot \frac{1}{(N-1)!} \; x^{N-1} \, e^{-x} \, dx \,, \tag{11}$$

when x stands for

$$\frac{N(1-2p)}{1-2\alpha} \tag{12}$$

with the understanding that when p becomes small the continuous Eulerian form is replaced by a terminating series with

$$p = {}^{r}\!/_{2N} \tag{13}$$

r taking integral values down to unity.

For this distribution it is easily found that

$$\left.\begin{aligned}
\Sigma(p) &= 1 \\
\Sigma p(1-2p) &= 1-2\alpha, & \Sigma(p^2) &= \alpha \\
\Sigma p(1-2p)^2 &= (1+\tfrac{1}{N})(1-2\alpha)^2, & \Sigma(p^3) &= \alpha^2 + \tfrac{1}{4N}(1-2\alpha)^2
\end{aligned}\right\} \tag{14}$$

Using this value in the estimated variance of an estimate of $\alpha$ derived from n observations, we have

$$\left.\begin{aligned}
V(\alpha) &= \frac{1}{n(n-1)} \left\{ 2s_2 + 4(n-2)s_3 - (4n-6)s_2^2 \right\} \\
&= \frac{1}{n(n-1)} \left\{ 2\alpha + 4(n-2)\alpha^2 + \frac{n-2}{N}(1-2\alpha)^2 \right. \\
&\qquad\qquad\qquad \left. - (4n-6)\alpha^2 \right\}
\end{aligned}\right\} \tag{15}$$

$$\text{---} \quad = \quad \frac{2\alpha(1-\alpha)}{n(n-1)} + \frac{n-2}{n(n-1)N}(1-2\alpha)^2 \qquad \Big) \quad (15)$$

or, for large populations, supposedly nearly in equilibrium

$$V(\alpha) = \frac{2\alpha(1-\alpha)}{n(n-1)} \quad ; \qquad (16)$$

the precision is therefore about that of a frequency based on $\frac{1}{2}n(n-1)$ observations.

For example, in 25 plants of a population of red clover, Bateman found 36 alleles occurring once only, while 7 occurred twice. Treating these as 50 independent observations

$$\sum a(a-1) = 14 \qquad (17)$$
$$\sum n(n-1) = 2450$$

giving the estimate of $\alpha$

$$\frac{1}{175} = \cdot 005714 28. \qquad (18)$$

The population size $N$ may be presumed to be large compared with

$$\frac{1}{2} \cdot 48 \cdot 735 = 3240, \qquad (19)$$

so the second term in the expression for the variance is ignored; the variance is then

$$V(\alpha) = \frac{\alpha(1-\alpha)}{1225} = \cdot 0^5 463807 \qquad (20)$$

and the standard error of random sampling is

$$\cdot 002154$$

For 7 observations out of 1225, Stevens' Table of approximate fiducial limits, the equivalent binomial observation, gives 2.81 and 14.38,

or odds of 19 to 1 that $\alpha$ lies between the limits

$$\cdot 00229 \quad \text{and} \quad \cdot 01174.$$

The number of alleles should then lie between

$$85 \quad \text{and} \quad 436$$

It would have made little difference if, as might be preferred to allow for the fact that in a sample of heterozygotes the observed frequencies are not quite independent the estimate had been taken to be

$$14/50 \times 48$$

instead of

$$14/50 \times 49.$$