

FINDING THE SIZE OF A FINITE POPULATION

BY D. A. DARLING¹ AND HERBERT ROBBINS²

University of California, Berkeley

1. Introduction. There are fixed sample size methods for estimating the unknown size N of a finite population by tagging the elements of a first sample and then counting the number of tagged elements of a second sample [1, p. 43]. Less well known are sequential methods [2], which have the advantage that the total sample size automatically adjusts itself to the unknown N to assure a desired accuracy of the estimate. All these methods only provide estimates for which the *relative* error is likely to be small. Suppose, however, that we want P_N (estimate = N) $\geq \alpha = .99$, say, no matter what the value $N = 1, 2, \dots$. How can this be done?

If we take as our estimate of N the number of distinct elements actually observed, the problem is one of finding a stopping rule such that the probability of having observed all N elements by the time we stop is $\geq \alpha$ for all $N \geq 1$. A concept of asymptotic efficiency may be introduced by comparing as $N \rightarrow \infty$ the expected sample size for any such rule with the fixed sample size necessary to observe all N elements with probability α . We give a procedure which is asymptotically efficient in this sense. We do not discuss the problem of finding a procedure which minimizes the Bayes expectation of the sample size for a given prior distribution of N .

Before going on, the reader is invited to consider the following problem. Sample one element at a time with replacement, tagging each element observed so that it can be recognized if it appears again. Choose some large integer M , and stop sampling when for the first time a run of M consecutive tagged elements occurs; estimate N to be the total number of distinct elements observed. Can we choose M so large that P_N (estimate = N) $\geq .99$ for all $N = 1, 2, \dots$? (Answer at end of paper.)

2. A procedure based on individual waiting times. An urn contains N white balls ($N = 1, 2, \dots$). We repeatedly draw a ball at random, observe its color, and replace it by a black ball. Eventually all N white balls will have been drawn and the urn will contain only black balls. The probability that this will occur at or before the n th draw is [1, pp. 92-93]

$$(1) \quad P_{N,n} = \sum_{i=0}^N (-1)^i \binom{N}{i} (1 - i/N)^n,$$

and if $N, n \rightarrow \infty$ so that $N e^{-n/N} \rightarrow \lambda, 0 < \lambda < \infty$, then

$$(2) \quad P_{N,n} \rightarrow e^{-\lambda}.$$

For any fixed $0 < \alpha < 1$ the smallest $n = n(N, \alpha)$ such that $P_{N,n} \geq \alpha$ can be found from (1) by trial and error. For large N this is tedious, but it follows from

Received 7 December 1966.

¹ Supported in part by National Science Foundation Grant GP 6549.

² Visiting Miller Research Professor.

(2) that $n = N \log N + cN + o(N)$, where c is determined by the equation $e^{-e^{-c}} = \alpha$. Equivalently, if we denote by Y_N the number of the draw on which the last (N th) white ball is drawn, then for any $-\infty < c < \infty$,

$$(3) \quad P_N((Y_N - N \log N)/N \leq c) \rightarrow e^{-e^{-c}} \quad \text{as } N \rightarrow \infty.$$

Suppose now that N is unknown and that we wish to find a rule for deciding when to stop drawing such that for a given $0 < \alpha < 1$,

$$(4) \quad P_N(\text{all } N \text{ white balls drawn by the time we stop}) \geq \alpha \quad (N \geq 1).$$

No fixed sample size will do this, so we must look for a sequential procedure.

Let Y_n denote the number of the draw on which the n th white ball appears ($n = 1, \dots, N$); thus $Y_1 \equiv 1$, and we put $Y_0 \equiv 0$, $Y_{N+1} \equiv \infty$ by convention. Define the waiting times

$$(5) \quad X_n = Y_{n+1} - Y_n \quad (n = 0, \dots, N),$$

so $X_0 \equiv 1$ and $X_N \equiv \infty$. The random variables X_1, \dots, X_{N-1} are independent, with the geometric distributions

$$(6) \quad P_N(X_n > j) = (n/N)^j \quad (j = 0, 1, \dots),$$

and

$$(7) \quad E_N(X_n) = N/(N - n), \quad \text{Var}_N(X_n) = nN/(N - n)^2.$$

Let (b_n) be any sequence of positive integers, and define B_n to be the event $X_n > b_n$ ($n = 1, \dots, N$). Since $X_N \equiv \infty$, B_N is certain. Let $J = \text{first } n \geq 1 \text{ such that } B_n \text{ occurs}$; then $1 \leq J \leq N$. Suppose we agree to stop drawing as soon as B_J occurs. All N white balls will have been drawn by the time we stop if and only if $J = N$. Hence the left hand side of (4) equals

$$(8) \quad P_N(J = N) = P_N(\cap_{n=1}^{N-1} B_n) = \prod_{n=1}^{N-1} P(X_n \leq b_n) = \prod_{n=1}^{N-1} \{1 - (n/N)^{b_n}\}.$$

It is clear that we can satisfy (4) by choosing the b_n properly. One way is the following. Define $b_1^* = \text{smallest integer } b_1 \text{ such that}$

$$1 - (\tfrac{1}{2})^{b_1} \geq \alpha.$$

Then (4) holds for $N = 2$ (and for $N = 1$ no matter what the sequence (b_n)). If b_1^*, \dots, b_{j-1}^* have been defined, set $b_j^* = \text{smallest integer } b_j \text{ such that}$

$$\prod_{n=1}^{j-1} \{1 - (n/(j+1))^{b_n^*}\} \cdot \{1 - (j/(j+1))^{b_j}\} \geq \alpha.$$

The sequence (b_n^*) gives a "step-wise minimal" solution of the inequalities

$$(9) \quad P_N(J = N) = \prod_{n=1}^{N-1} \{1 - (n/N)^{b_n}\} \geq \alpha \quad (N = 2, 3, \dots).$$

(It is not *uniformly* minimal, since if (b_n) satisfies (9) then we can increase b_1, \dots, b_{j-1} sufficiently so that a smaller b_j will still satisfy (9).) Although the b_n^* are not given by an explicit formula they can be computed numerically for any given α .

Instead of doing this we shall find a lower bound for the left hand side of (9)

for the particular sequence

$$(10) \quad b_n = [cn] + 1 \quad (n = 1, 2, \dots),$$

where c is a suitable positive constant. We shall show that if c is chosen large enough then (9) holds, and shall find the limit of $P_N(J = N)$ as $N \rightarrow \infty$.

For the sequence (10) let β be any constant $0 < \beta < 1$. Then

$$(11) \quad P_N(J = N) = \prod_{1 \leq n \leq \beta N} \{1 - (n/N)^{b_n}\} \cdot \prod_{\beta N < n \leq N-1} \{1 - (n/N)^{b_n}\} = Q_N \cdot R_N,$$

where

$$(12) \quad Q_N \geq \prod_{1 \leq n \leq \beta N} (1 - \beta^{cn}) \geq \prod_{n=1}^{\infty} (1 - \beta^{cn}),$$

and since $\log(1 - x) \leq -x$,

$$(13) \quad \begin{aligned} R_N &= \prod_{1 \leq i < (1-\beta)N} \{1 - (1 - i/N)^{[c(N-i)]+1}\} \\ &\geq \prod_{1 \leq i < (1-\beta)N} \{1 - (1 - i/N)^{c(N-i)}\} \\ &= \prod_{1 \leq i < (1-\beta)N} \{1 - e^{c(N-i) \log(1-i/N)}\} \\ &\geq \prod_{1 \leq i < (1-\beta)N} \{1 - e^{-ci(1-i/N)}\} \\ &\geq \prod_{1 \leq i < (1-\beta)N} \{1 - e^{-ci\beta}\} \geq \prod_{i=1}^{\infty} (1 - e^{-ci\beta}). \end{aligned}$$

Defining for $0 \leq x < 1$ the function

$$(14) \quad \varphi(x) = \prod_{i=1}^{\infty} (1 - x^i) \geq 1 - \sum_{i=1}^{\infty} x^i = 1 - x/(1 - x) \rightarrow 1 \quad \text{as } x \rightarrow 0,$$

we have the uniform lower bound

$$(15) \quad P_N(J = N) \geq \varphi(\beta^c) \cdot \varphi(e^{-c\beta}) \quad (N = 2, 3, \dots).$$

In particular, if we choose β to be the root $\beta_0 = 0.56 \dots$ of the equation

$$(16) \quad \beta = e^{-\beta},$$

then

$$(17) \quad P_N(J = N) \geq \varphi^2(e^{-c\beta_0}).$$

Choosing c to satisfy $\varphi^2(e^{-c\beta_0}) = \alpha$ it follows that (9) holds.

We can improve (15) somewhat for large N . Write (13) in the form

$$(18) \quad \log R_N = \sum_{i=1}^{\infty} a_{i,N},$$

where

$$(19) \quad \begin{aligned} a_{i,N} &= \log \{1 - (1 - i/N)^{[c(N-i)]+1}\} \quad \text{for } 1 \leq i < (1 - \beta)N, \\ &= 0 \quad \text{for } i \geq (1 - \beta)N. \end{aligned}$$

For any fixed $i = 1, 2, \dots$

$$(20) \quad \lim_{N \rightarrow \infty} a_{i,N} = \log(1 - e^{-ci}) = a_i, \text{ say,}$$

and for $1 \leq i < (1 - \beta)N$ we have as in (13),

$$(21) \quad 0 \leq a_{i,N} \leq \log \{1 - (1 - i/N)^{c(N-i)}\} \\ \geq \log (i - e^{-ci(1-i/N)}) \geq \log (1 - e^{-ci\beta}),$$

so this holds for *all* i, N , and

$$(22) \quad \sum_{i=1}^{\infty} \log (1 - e^{-ci\beta}) = \log \varphi(e^{-c\beta}) > -\infty.$$

By the dominated convergence theorem,

$$(23) \quad \lim_{N \rightarrow \infty} R_N = \lim_{N \rightarrow \infty} \exp [\sum_{i=1}^{\infty} a_{i,N}] = \exp [\sum_{i=1}^{\infty} a_i] \\ = \prod_{i=1}^{\infty} (1 - e^{-ci}) = \varphi(e^{-c}).$$

Since by (12), $\varphi(\beta^c)R_N \leq P_N(J = N) \leq R_N$, it follows that

$$\varphi(\beta^c)\varphi(e^{-c}) \leq \liminf_{N \rightarrow \infty} P_N(J = N) \leq \limsup_{N \rightarrow \infty} P_N(J = N) \leq \varphi(e^{-c}),$$

and since β can be arbitrarily near 0, and $\varphi(\beta^c) \rightarrow 1$ as $\beta \rightarrow 0$,

$$(24) \quad \lim_{N \rightarrow \infty} P_N(J = N) = \varphi(e^{-c}).$$

Thus from (17),

$$(25) \quad \varphi^2(e^{-c\beta_0}) \leq P_N(J = N) \rightarrow \varphi(e^{-c}) \quad \text{as } N \rightarrow \infty.$$

For any $\epsilon > 0$, if we increase the first $j = j(\epsilon)$ terms of (10) we can clearly strengthen (25) to read

$$(26) \quad \varphi(e^{-c}) - \epsilon \leq P_N(J = N) \rightarrow \varphi(e^{-c}) \quad \text{as } N \rightarrow \infty.$$

To see how efficient this procedure is, let us look at its sample size

$$(27) \quad S = X_0 + \cdots + X_{J-1} + b_J \leq X_0 + \cdots + X_{N-1} + b_N.$$

From (7),

$$(28) \quad E_N(S) \leq N(1 + \tfrac{1}{2} + \cdots + 1/N) + cN + 1,$$

which is somewhat greater than the fixed sample size $n = [N \log N + cN]$ for which we have seen that $P_{N,n} \rightarrow e^{-e^{-c}}$. Now

$$0 < \varphi(e^{-c}) < 1 - e^{-c} < e^{-e^{-c}} < 1,$$

so (26) shows that for $N \rightarrow \infty$ the probability of having drawn all the white balls by the sequential procedure is somewhat less than the corresponding probability for the fixed sample size $n = [N \log N + cN]$. Of course, the latter procedure requires a knowledge of N ; moreover, the ratio of the two error probabilities is small for large values of c ;

$$(29) \quad (1 - e^{-e^{-c}})/(1 - \varphi(e^{-c})) \rightarrow 1 \quad \text{as } c \rightarrow \infty.$$

Nevertheless, the fact remains that the sequential procedure is somewhat inefficient for any fixed c and large values of N . We therefore ask, is there any

sequential procedure such that for fixed $-\infty < c < \infty$,

$$(30) \quad \text{Sample size} \leq N \log N + cN,$$

$$P_N(\text{all } N \text{ white balls drawn by the time we stop}) \rightarrow e^{-e^{-c}} \text{ as } N \rightarrow \infty?$$

An affirmative answer is given in the next section.

3. An asymptotically efficient procedure based on cumulative waiting times.

We modify the sequential procedure of the previous section as follows. Let (a_n) be a sequence of positive constants, and define A_n to be the event that $Y_{n+1} > a_n$ ($n = 1, \dots, N$). Since $Y_{N+1} \equiv \infty$, A_N is certain. Let $I = \text{first } n \geq 1 \text{ such that } A_n \text{ occurs}$. Then $1 \leq I \leq N$. We agree to stop as soon as A_I occurs. Then (cf. (8))

$$\begin{aligned} P_N(\text{all } N \text{ white balls drawn by the time we stop}) &= P_N(I = N) \\ (31) \quad &= P_N(\bigcap_{n=1}^{N-1} A_n') = P_N(\bigcap_{n=1}^{N-1} (Y_{n+1} \leq a_n)) \\ &= P_N(\bigcap_{n=1}^{N-1} (X_0 + \dots + X_n \leq a_n)). \end{aligned}$$

As before, we could define a "step-wise minimal" sequence (a_n^*) such that the expression (31) is $\geq \alpha$ for every $N \geq 2$, but since the events A_n are not independent the explicit computation of the a_n^* is difficult. Instead, as before, we shall choose an explicit sequence (a_n) and estimate the value of (31) when N is large. Our sequence is the following. Let c be any finite constant, let $n^* = \text{smallest } n \text{ such that } n \geq e^{1-c}$, and define

$$\begin{aligned} (32) \quad a_n &= n + 1 \quad \text{for } n = 1, \dots, n^* - 1, \\ &= n \log n + cn \quad \text{for } n \geq n^*. \end{aligned}$$

It is clear that for this procedure the sample size is always $\leq a_N (= N \log N + cN \text{ for } N \geq n^*)$. And we shall prove that $P_N(I = N) \rightarrow e^{-e^{-c}}$ as $N \rightarrow \infty$. (It will be seen from the proof that as in (26) for any $\epsilon > 0$ we can increase the first $j = j(\epsilon)$ values of (a_n) so as to make

$$e^{-e^{-c}} - \epsilon \leq P_N(I = N) \rightarrow e^{-e^{-c}} \text{ as } N \rightarrow \infty.)$$

It is easy to check that

$$\begin{aligned} (33) \quad a_n - a_{n-1} &\geq 1 \quad \text{for all } n \geq 2, \\ &\geq \log n + c \quad \text{for } n \geq n^*, \end{aligned}$$

and that the random variables

$$(34) \quad w_n = X_{N-n}/N - 1/n \quad (n = 1, \dots, N)$$

are independent with

$$(35) \quad E_N(w_n) = 0, \quad \text{Var}_N(w_n) = 1/n^2 - 1/Nn < 1/n^2.$$

By Kolmogorov's inequality, for any $d > 0$,

$$(36) \quad P_N(w_1 + \cdots + w_n \leq -d \text{ for some } n = 1, \cdots, N) \\ \leq (\sum_1^\infty 1/n^2)/d^2 = \pi^2/6d^2.$$

We have

$$(37) \quad A_n' = (X_0 + \cdots + X_n \leq a_n) = (w_{N-n} + \cdots + w_N \leq b_n), \\ b_n = a_n/N - ((N-n)^{-1} + \cdots + 1/N).$$

And for $n \geq 2$,

$$(38) \quad A_{N-1}' \cap (w_1 + \cdots + w_n \geq b_{N-1} - b_{N-n}) \subset A_{N-n}'.$$

Now for $N \geq n + n^*$,

$$(39) \quad b_{N-1} - b_{N-n} = ((a_{N-1} - a_{N-n})/N) + (1 + \frac{1}{2} + \cdots + 1/(n-1)) \\ = (n-1) \log N/N - (1 + \frac{1}{2} + \cdots + 1/(n-1)) \\ + c(n-1)/N + f(1 - n/N) - f(1 - 1/N),$$

where

$$(40) \quad 0 < f(x) = -x \log x < e^{-1} \quad \text{for } 0 < x < 1.$$

Hence for $N \geq n + n^*$,

$$(41) \quad b_{N-1} - b_{N-n} \leq n \log N/N - \log n + \beta \quad (\beta = |c| + e^{-1}).$$

Put

$$(42) \quad p = d + \beta + 1, \quad k = e^p.$$

Then as $N \rightarrow \infty$,

$$k \log N/N - \log k + \beta \rightarrow -d - 1,$$

$$(N(1 - p/\log N)/N) \log N - \log(N(1 - p/\log N)) + \beta \rightarrow -d - 1,$$

and for $N \geq N_d$,

$$(43) \quad b_{N-1} - b_{N-n} \leq -d \quad \text{for } k \leq n \leq N(1 - p/\log N).$$

Hence from (38) and (36),

$$(44) \quad P_N(A_{N-n}' \text{ for all } 1 \leq n \leq N(1 - p/\log N)) \\ \geq P_N(A_{N-1}' \cap \cdots \cap A_{N-k}') \\ - P_N(w_1 + \cdots + w_n \leq -d \text{ for some } n = 1, \cdots, N) \\ \geq P_N(A_{N-1}' \cap \cdots \cap A_{N-k}') - \pi^2/6d^2.$$

But from (3), as $N \rightarrow \infty$

$$(45) \quad P_N(A_{N-1}' \cap \cdots \cap A_{N-k}') \geq P_N(Y_N \leq a_{N-k}) \\ = P_N((Y_N - N \log N)/N \leq ((a_{N-k} - N \log N)/N) \rightarrow e^{-e^{-c}},$$

since as $N \rightarrow \infty$

$$\begin{aligned} (a_{N-k} - N \log N)/N &= ((N-k) \log (N-k) - c(N-k) - N \log N)/N \\ &= (1 - k/N) \log (1 - k/N) + (1 - k/N) \log N \\ &\quad - \log N - c(1 - k/N) \rightarrow -c. \end{aligned}$$

Hence

$$(46) \quad \liminf_{N \rightarrow \infty} P_N(A'_N \text{ for all } 1 \leq n \leq N(1 - p/\log N)) \geq e^{-e^{-c}} - \pi^2/6d^2.$$

We shall show in a moment that

$$(47) \quad \lim_{N \rightarrow \infty} P_N(A'_n \text{ for all } 1 \leq n \leq pN/\log N) = 1.$$

It will then follow from (44) and (45) that

$$(48) \quad \liminf_{N \rightarrow \infty} P_N(A'_n \text{ for all } 1 \leq n \leq N-1) \geq e^{-e^{-c}} - \pi^2/6d^2.$$

Since d can be arbitrarily large, (48) holds without the last term. But

$$(49) \quad P_N(A'_n \text{ for all } 1 \leq n \leq N-1) \leq P_N(A'_{N-1}) \rightarrow e^{-e^{-c}}$$

by (45) for $k = 1$. Hence

$$(50) \quad \lim_{N \rightarrow \infty} P_N(A'_n \text{ for all } 1 \leq n \leq N-1) = \lim_{N \rightarrow \infty} P_N(I = N) = e^{-e^{-c}}.$$

It remains only to prove (47). Now setting $a_0 = 1$

$$\begin{aligned} P_N(A'_n \text{ for all } 1 \leq n \leq pN/\log N) &= P_N(\bigcap_{1 \leq n \leq pN/\log N} (X_0 + \cdots + X_n \leq a_n)) \\ &\geq P_N(\bigcap_{1 \leq n \leq (pN/\log N)} (X_n \leq a_n - a_{n-1})) \\ &= \prod_{1 \leq n \leq pN/\log N} \{1 - (n/N)^{a_n - a_{n-1}}\} \\ &\geq 1 - \sum_{1 \leq n \leq pN/\log N} (n/N)^{a_n - a_{n-1}}, \end{aligned}$$

and by (33), as $N \rightarrow \infty$

$$\begin{aligned} \sum_{1 \leq n \leq N^{1/2}/\log N} (n/N)^{a_n - a_{n-1}} &\leq \sum_{1 \leq n \leq N^{1/2}/\log N} (n/N) \leq 1/(\log N)^2 \rightarrow 0, \\ \sum_{N^{1/2}/\log N < n \leq pN/\log N} (n/N)^{a_n - a_{n-1}} &\leq \sum_{N^{1/2}/\log N < n \leq pN/\log N} (p/\log N)^{\log n + c} \\ &\leq (p/\log N)^c \sum_{n > N^{1/2}/\log N} n^{\log p - \log \log N} \\ &\leq (p/\log N)^c \cdot \int_{N^{1/2}}^{\infty} (dx/x^2) \rightarrow 0, \end{aligned}$$

which completes the proof of (47).

The answer to the question in Section 1 is no; for any M ,

$$\lim_{N \rightarrow \infty} P_N(\text{estimate} = N) = 0.$$

REFERENCES

- [1] FELLER, W. (1957). *Introduction to Probability Theory and its Applications*. **1** (2nd ed.) Wiley, N. Y.
- [2] GOODMAN, L. A. (1953). Sequential Sampling Tagging for Population Size Problems. *Ann. Math. Statist.* 56-69.