

This article presents a procedure and a table for selecting sample size for simultaneously estimating the parameters of a multinomial distribution. The results are obtained by examining the "worst" possible value of a multinomial parameter vector, analogous to the case in which a binomial parameter equals one-half.

KEY WORDS: Multinomial distribution; Simultaneous inference.

1. INTRODUCTION

Sample size problems rarely have satisfyingly simple answers, because the anticipated precision of estimates may depend on one or more unknown parameters. An exception in which a definite answer can be given in the absence of prior information is the binomial distribution, for which we can compute the sample size sufficient to achieve a specified degree of precision for the "worst case" parameter value of one-half. For the multinomial distribution, with more than two categories, it is not realistic to consider a "worst case" in which all parameters equal one-half, since the sum of the parameters must equal one. In this article I answer the question, what is the "worst" possible value of a multinomial parameter vector, and give a table of sample sizes specified precision levels.

The problem of choosing a sample size for the simultaneous estimation of multinomial proportions is closely akin to the problem of simultaneous confidence intervals for multinomial proportions, with the difference that in the sample size problem the acceptable widths of the intervals are specified in advance and sample size is chosen to control the probability that the intervals will cover the true proportions. Queensbury and Hurst (1964) presented a method for constructing simultaneous confidence intervals for multinomial proportions based on the approximate chi-squared distribution of the sum of the observed minus expected frequencies squared divided by the expected frequencies. Goodman (1965) improved on this result by constructing less conservative (shorter) intervals satisfying the stated level of significance. Goodman's method was based on the normal approximation for a binomial proportion and used Bonferroni's inequality to put a bound on the probability that all of the intervals would be simultaneously correct. Goodman did not address the problem of sample size selection, but his approach has been the basis for subsequent work on the subject.

Angers (1974) applied Goodman's method to the problem of sample size for multinomial proportions and presented a

*Steven K. Thompson is Assistant Professor, Department of Mathematical Sciences, University of Alaska, Fairbanks, AK 99775-1110. The work for this article was initiated while the author was Biometrician, Alaska Department of Fish and Game, Kodiak, AK. The author would like to thank the referees for valuable suggestions.

graphical method for selecting a sample size based on prior knowledge of the parameter values. Tortora (1978) based sample size on the "worst case" individual parameter, which is the parameter closest to .5 when the precision criteria are the same for all parameters. Angers (1979) pointed out that Tortora's method was more conservative than necessary in some cases. For selecting sample size with limited prior information about the parameters in the form of inequalities, Angers proposed using the value closest to .5 for each parameter. Angers (1979, 1984) described the general procedure for selecting sample size using prior estimates of parameter values. The method is computationally tedious to carry out but results in substantial reductions in sample size over previous methods. For equal interval widths, Angers (1984) gave an empirical result on the "worst case" parameter vectors for small α levels, based on Monte Carlo selection of large numbers of parameter vectors.

In this article I establish the form of the "worst case" multinomial parameter vector when equal interval widths are specified for each parameter component. A formula for sample size under this worst case is given, and a table is provided that makes sample size easy to determine for selected significance levels and any interval width. Angers' empirical result is proved and extended to all α levels. The form of the worst case parameter vector is also established for the case in which the possible parameter values are constrained by prior knowledge to satisfy inequalities. It is noted that when different acceptable confidence interval widths are specified for different parameters, determination of sample size remains a computational problem. Examples from the literature are reworked, illustrating the possible reductions in sample size or simplification of the sample size determination procedure using the methods of this article.

2. METHOD

The objective is to select the smallest sample size n for a random sample from a multinomial population such that the probability will be at least $1 - \alpha$ that all of the estimated proportions will simultaneously be within specified distances of the true population proportions, that is,

$$\Pr\left\{\bigcap_{i=1}^k |p_i - \pi_i| \leq d_i \mid \pi_i \leq 1 - \alpha, -2 \leq p_i \leq 2\right\} \geq 1 - \alpha, \quad \text{where } \pi_i \text{ is the proportion in the } i\text{th category in the population, } p_i \text{ is the observed proportion, and } k \text{ is the number of categories.}$$

In this article, it is assumed that the population is large enough for finite population correction factors to be ignored if sampling is done without replacement and that sample sizes are large enough for the normal approximation to be used. For an individual parameter π_i , the probability α_i that the estimate p_i lies outside the specified interval is, by the nor-

*James R. Waters and Alexander J. Chester are with the National Marine Fisheries Service, Southeast Fisheries Center, Beaufort Laboratory, Beaufort, NC 28516. The authors thank David Colby, Douglas Vaughan, and two anonymous reviewers for their helpful comments on the manuscript.

Our purpose is to determine the optimal allocation of sampling effort in a two-stage design when several parameters are estimated simultaneously. We begin with a brief discussion of optimal survey design when estimating a single

Analyses of two-stage designs are well documented when a single variable is measured, and methods to calculate the optimal allocation of sampling effort among sampling stages are readily available (Cochran 1977, chap. 10; Snedecor and Cochran 1980, chap. 21; Sokal and Rohlf 1981, chap. 10). Many surveys, however, measure several variables simultaneously, and procedures for determining optimal sampling effort are not well defined. The traditional approach has been to estimate optimal sample size for each variable individually and then choose the final sampling design from among the individual solutions. The shortcoming of this approach is that it does not necessarily identify all possible candidates for the overall solution.

1. INTRODUCTION

Optimal Allocation in Multivariate, Two-Stage Sampling Designs

JAMES R. WATERS and ALEXANDER J. CHESTER*

An inequality involving Mills ratio (Patel and Read 1982,

$$g'(u) = -2au(1-u)\phi^{-1} + [\phi^{-1}(1-au)]^2(1-2u).$$

$$\times (1-au)[\phi^{-1}(1-au)]^{-1}$$

$$= 1/m \text{ and } a = \alpha/2.$$

 Proof. Let $g(u) = [\phi^{-1}(1-au)]^2u(1-u)$, where u

$$[\phi^{-1}(1-\alpha/2m)]^2 \geq 2(1-1/m)/(1-2/m).$$

The function $f(m) = \{\phi^{-1}(1-\alpha/2m)\}^2(1/m)(1-1/m)$ decreases for all $m > 2$ such that

APPENDIX B: A RESULT USEFUL IN SAMPLE SIZE COMPUTATION

remaining π_i on the boundary of A .
 vector with $\pi_i = (1-b)/m$ for m parameters and the interior of A , an extremum is provided by the parameter boundary region of A under consideration. If π^* is in the not in the interior of A , there is no extremum point in the equal, that is, $\pi_i = (1-b)/m$ ($i = 1, \dots, m$). If π^* is extremum point is provided by π^* , in which the π_i are same as in the proof of Theorem 1, and the only interior strain $\Sigma \pi_i = b$, the derivative equations of G are exactly Set $G(\pi, \lambda) = \Sigma \alpha_i - \lambda(\Sigma \pi_i - b)$. Except for the $\Sigma \pi_i = b$, for the m remaining parameters, subject to $\Sigma \pi_i = 1 - b$.
 problem of maximizing $\Sigma \pi_i \alpha_i$ becomes that of maximizing α_i is fixed. Denote by b the sum of these $k - m$ α_i . The

Angers, Claude (1974). "A Graphical Method to Evaluate Sample Sizes for the Multinomial Distribution," *Technometrics*, 16, 469-471.
 (letter to the editor), *The American Statistician*, 33, 163-164.
 (1984). "Large Sample Sizes for the Estimation of Multinomial Frequencies From Simulation Studies," *Simulation*, October, 175-178.
 Cochran, William G. (1977). *Sampling Techniques* (3rd ed.). New York: John Wiley.
 Goodman, Leo A. (1965). "On Simultaneous Confidence Intervals for Multinomial Populations," *Technometrics*, 7, 247-254.
 Hancok, Harris (1960). *Theory of Maxima and Minima*. New York: Dover Publications.
 Miller, Rupert G., Jr. (1980). *Simultaneous Statistical Inference* (2nd ed.). New York: Springer-Verlag.
 Patel, J. K., and Read, C. B. (1982). *Handbook of the Normal Distribution*. New York: Marcel Dekker.
 Queensbury, C. P., and Hurst, D. C. (1964). "Large Sample Simultaneous Confidence Intervals for Multinomial Proportions," *Technometrics*, 6, 191-195.
 Torner, Robert D. (1978). "A Note on Sample Size Estimation for Multinomial Populations," *The American Statistician*, 32, 100-102.

REFERENCES

p. 64) gives $\pi_i \phi(x) \leq 1/(1 - \phi(x))$ for $x > 0$, so with $x = \phi^{-1}(1-au)$,

$$g(u) \geq -2(1-u) + [\phi^{-1}(1-au)]^2(1-2u).$$

 Thus g is increasing and f is decreasing whenever the right side is positive, that is, for all $m > 2$ such that

$$[\phi^{-1}(1-\alpha/2m)]^2 \geq 2(1-1/m)/(1-2/m).$$

 [Received July 1984. Revised December 1985.]