

BOUNDED LENGTH CONFIDENCE INTERVALS FOR THE p -POINT OF A DISTRIBUTION FUNCTION, II¹

BY R. H. FARRELL

Cornell University

1. Introduction. Let $0 < p < 1$. A number $\gamma_{p,F}$ is a p -point of the distribution function F if $F(\gamma_{p,F}) \geq p$ while $F(\gamma_{p,F}-) \leq p$. Given $L > 0$ and $0 < \alpha < 1$, a (sequential) confidence interval procedure is a L - α bounded length confidence interval procedure if when sampling stops an interval of length not exceeding L is given which covers $\gamma_{p,F}$ with probability at least $1 - \alpha$.

Throughout this paper we consider only those decision procedures that are based on independently and identically distributed random variables. In the sequel we make this completely precise. We present a negative result (see the statement of the theorem below). It will be convenient to say simply "there does not exist a procedure that works for all $F \in \mathbf{F}$." This means a set \mathbf{F} of distribution functions is specified and that the common distribution function F of the random variables is in \mathbf{F} . The problem is to obtain a confidence interval of length $\leq L$ for the p -point $\gamma_{p,F}$ of the common distribution function F . The experimenter is allowed to construct his procedure using the information $F \in \mathbf{F}$. We show that there cannot exist confidence interval procedures of a specified type giving an interval of length $\leq L$ yet covering $\gamma_{p,F}$ with probability $\geq 1 - \alpha$ for all $F \in \mathbf{F}$. Thus, for example, it is clear that no sequential confidence interval procedure can work for all F satisfying $F(\gamma_{p,F} - L-) = F(\gamma_{p,F} + L)$.

In Section 2 of the first of this series of papers, Farrell [3], we define a measure of flatness by

$$(1.1) \quad \epsilon_F = \sup_{0 < \rho < 1} \min (F(\gamma_{p,F} + \rho L) - p, p - F(\gamma_{p,F} + (\rho - 1)L)).$$

We shall be interested in confidence interval procedures that may be applied to observations on $F \in \mathbf{F}$, where if $F \in \mathbf{F}$ then $\epsilon_F > 0$.

We are interested in choices of \mathbf{F} for which fixed sample size procedures will fail to be L - α bounded length confidence interval procedures. In case $F \in \mathbf{F}$ implies that F has a unimodal density function, Weiss [5] has shown the existence of two-stage L - α bounded length confidence interval procedures. In the general case of $F \in \mathbf{F}$ if and only if $\epsilon_F > 0$ we show in [4] the construction of a sequential L - α bounded length confidence interval procedure with certain optimality properties. Examples of sequential procedures have also been constructed by J. Kiefer and L. Weiss but these examples have not been published.

The present paper gives a nonexistence result.

Received 3 September 1965.

¹ Research sponsored in part by the Office of Naval Research under Contract Nonr 401(50) with Cornell University. Part of this work appears in the author's Ph.D. dissertation, University of Illinois.

THEOREM. *If $m \geq 1$ and if \mathbf{F} contains all distribution functions F having bimodal density functions then there does not exist an m -stage L - α bounded length confidence interval procedure that works for all $F \in \mathbf{F}$.*

A similar result has been proven by J. Kiefer and L. Weiss but their result has not been published. The author [3], Theorem 2, has a related result. The author wishes to thank D. L. Burkholder for suggesting this problem and for helpful conversations.

2. Proof of the theorem. An m -stage procedure consists of the sampling plan (stopping rule) and a function Δ of the observations which is the terminal decision. In this section we develop a series of inequalities. The last of these inequalities shows that given an m -stage sampling plan, given a function Δ of the observations, and given an integer $n \geq 1$, there is an $F \in \mathbf{F}$ such that $P(|\Delta - \gamma_{p,F}| < L/2) < 1/n$.

We begin by giving a precise description of the terminology. Suppose $\{X_n, n \geq 1\}$ is a sequence of independently and identically distributed random variables each having F as distribution function. Suppose random variables N and Y are given such that $N \geq 1$ and N is integer valued. Suppose

(1) the conditional distribution of N and Y given $\{X_n, n \geq 1\}$ exists and does not depend on F ;

(2) for all real numbers y and integers i the event $\{N = i, Y \leq y\}$ is independent of the collection of random variables $\{X_n, n \geq i + 1\}$. Then we shall say that Y is determined by a sequential sampling plan. For future reference note that if a random variable λ is defined by $\lambda = 1$ when $N \leq k$ and $\lambda = 0$ when $N > k$, then the joint conditional distributions of λ, Y given $\{X_n, n \geq 1\}$ does not depend on F . Therefore $E(\lambda Y | X_1, X_2, \dots) = E(\lambda Y | X_1, \dots, X_k)$ does not depend on F .

The definition of an m -stage sampling plan is very similar. Let N_1, \dots, N_m be integer valued random variables such that for $1 \leq i \leq m, N_i \geq 1$. Let $N = N_1 + \dots + N_m$. Suppose Y is a random variable such that

(1) The conditional distribution of N_1, \dots, N_m , given $\{X_n, n \geq 1\}$ is independent of F ;

(2) N_1 is independent of $\{X_n, n \geq 1\}$;

(3) The event $N_1 + \dots + N_k = i, N_{k+1} = j$ is independent of X_{i+1}, X_{i+2}, \dots for $i \geq 1, j \geq 1, 1 \leq k \leq m - 1$;

(4) For all real y and integers $i \geq 1$ the event $N = i$ and $Y \leq y$ is independent of X_{i+1}, X_{i+2}, \dots .

Then we shall say that Y is determined by a m -stage sampling plan. It should be noted that Y is determined by a sequential sampling plan. Also, notice that N is determined by an $(m - 1)$ -stage sampling plan.

These definitions of sequential and m -stage sampling are broad enough to include randomized sampling and decision rules.

In the following F and G will denote distribution functions and S will be the set of real numbers x such that $F(x) \neq G(x)$. R will denote the set of real

numbers and E_k will denote Euclidean k -space. The probabilities of certain events will be compared when the random variables $\{X_n, n \geq 1\}$ each have F or each have G as distribution. Probabilities and expectations relative to these distributions will be indicated by using “ F ” or “ G ” as a subscript.

LEMMA 2.1. *Suppose S is contained in an interval (a, b) whose endpoints are points of continuity of F and G . Suppose φ is a real valued Borel measurable function on E_k and that $0 \leq \varphi \leq 1$. Further suppose $F(b) - F(a) \leq \epsilon \leq 1$ and $G(b) - G(a) \leq \epsilon$. Then $|E_F(\varphi) - E_G(\varphi)| \leq 1 - (1 - \epsilon)^k \leq k\epsilon$.*

PROOF. Let $W = R - (a, b)$ and $T = W \times W \times \cdots \times W$, the Cartesian product k times. Observe for use in the following that if ψ is the characteristic function of the set T then $E_F(\varphi\psi) = E_G(\varphi\psi)$ since F and G induce the same k -dimensional measures on the Borel subsets of T . Then

$$\begin{aligned} E_F(\varphi) &= E_F(\varphi\psi) + E_F(\varphi(1 - \psi)) \\ &= E_G(\varphi\psi) + E_F(\varphi(1 - \psi)) \\ &\leq E_G(\varphi) + P_F(\psi = 0) \\ &= E_G(\varphi) + 1 - P_F(X_1 \in W) \cdots P_F(X_k \in W) \\ &\leq E_G(\varphi) + 1 - (1 - \epsilon)^k. \end{aligned}$$

By interchanging F and G in the above inequality the new inequality

$$E_G(\varphi) \leq E_F(\varphi) + 1 - (1 - \epsilon)^k$$

is obtained. These two inequalities are equivalent to the first inequality of the lemma. Since $0 \leq \epsilon \leq 1$, the inequality $1 - (1 - \epsilon)^k \leq k\epsilon$ follows at once from the mean value theorem.

LEMMA 2.2. *Let $F, G, S, (a, b)$ and ϵ be as in Lemma 2.1. Let the random variable φ be determined by a sequential sampling plan and suppose $0 \leq \varphi \leq 1$. Then for each integer $k \geq 1$, $|E_F(\varphi) - E_G(\varphi)| \leq 1 - P_F(N \leq k) + k\epsilon$.*

PROOF. The following method of proof and the notations used derive from the article of Bahadur and Savage [1]. Let $F^{(k)}$ denote the distribution on E_k corresponding to the coordinate variables being independent and distributed as F . Similarly for $G^{(k)}$. Define a random variable λ by $\lambda = 1$ when $N \leq k$ and $\lambda = 0$ when $N > k$. Then

$$E_G(\varphi) \geq E_G(\lambda\varphi) = \int_{E_k} E_G(\lambda\varphi | X_1, \dots, X_k) dG^{(k)}$$

As observed above $E(\lambda\varphi | X_1, \dots, X_k)$ does not depend on the distributions F or G . Therefore

$$\begin{aligned} E_G(\lambda\varphi) &= \int_{E_k} E_F(\lambda\varphi | X_1, \dots, X_k) dG^{(k)} \\ &\geq E_F(\lambda\varphi) - k\epsilon \quad (\text{by the previous lemma}) \\ &= E_F(\varphi) - E_F(\varphi - \lambda\varphi) - k\epsilon \\ &\geq E_F(\varphi) - E_F(1 - \lambda) - k\epsilon. \end{aligned}$$

Therefore $E_F(\varphi) - E_G(\varphi) \leq 1 - P_F(N \leq k) + k\epsilon$. Replace φ by $1 - \varphi$. The above arguments remain valid and the new inequality $E_G(\varphi) - E_F(\varphi) \leq 1 - P_F(N \leq k) + k\epsilon$ is obtained. Thus the lemma is proven.

LEMMA 2.3. Suppose $\delta > 0$. Let the random variable φ be determined by a one-stage sampling plan and suppose $0 \leq \varphi \leq 1$. Then there is an $\epsilon > 0$ such that for every pair of distributions F and G and every pair of real numbers a, b if

- (1) $F(b) - F(a) \leq \epsilon$;
- (2) $\{x \mid F(x) \neq G(x)\}$ is contained in (a, b) ;
- (3) the points a and b are continuity points of F and G ;

then $|E_F(\varphi) - E_G(\varphi)| \leq \delta$.

PROOF. The distribution of N does not depend on F and G . Choose an integer $k \geq 1$ such that $P(N \leq k) \geq 1 - \delta/2$. Choose $\epsilon > 0$ so small that $\epsilon k < \delta/2$. Then Lemma 2.3 follows from Lemma 2.2.

LEMMA 2.4. Let the random variable φ be determined by an m -stage sampling plan and let $0 \leq \varphi \leq 1$. Suppose $\delta > 0$ and $0 < p < 1$. Let a real number interval (a, b) be given. There exists a distribution function F with bimodal density function f having modes α and β satisfying, $\alpha < a < b < \beta$, $F(\alpha) < p < F(\beta)$, and $F(\beta) - F(\alpha) > 0$. In addition, if G is any distribution function satisfying $\{x \mid F(x) \neq G(x)\} \subset (\alpha, \beta)$ and α and β are points of continuity of G then $|E_F(\varphi) - E_G(\varphi)| < \delta$.

PROOF. By induction on the number of stages m . The case $m = 1$ follows at once from Lemma 2.3. Assume the conclusion of Lemma 2.4 holds for $m - 1$. Relative to modes α^* , β^* and the number $(m - 1)\delta/m$, such that $\alpha^* < a < b < \beta^*$, let H satisfy the conclusion of the lemma.

The total sample size N is determined by an $(m - 1)$ -stage sampling plan. By the inductive hypothesis, if $k \geq 1$ then $|P_H(N \leq k) - P_F(N \leq k)| < (m - 1)\delta/m$ whenever $\{x \mid H(x) \neq F(x)\} \subset (\alpha^*, \beta^*)$ and α^* and β^* are continuity points of F . Choose k^* so that $1 - P_H(N \leq k^*) \leq \delta/(2m)$. Choose $\epsilon > 0$ so that $\epsilon k^* < \delta/(2m)$ and $\epsilon \leq H(\beta^*) - H(\alpha^*)$. In addition take ϵ so small that we may find a distribution function F having a bimodal density function with modes α and β satisfying $\alpha^* \leq \alpha < a < b < \beta \leq \beta^*$, $F(\beta) - F(\alpha) = \epsilon$, and $F(\alpha) < p < F(\beta)$. We may obtain F by modification of H , so that $\{x \mid H(x) \neq F(x)\} \subset (\alpha^*, \beta^*)$. Then for any distribution function G satisfying $\{x \mid F(x) \neq G(x)\} \subset (\alpha, \beta)$ and for which α and β are continuity points of G , we obtain from Lemma 2.2, $|E_F(\varphi) - E_G(\varphi)| \leq 1 - P_F(N \leq k^*) + \epsilon k^* < 1 - P_H(N \leq k^*) + (m - 1)\delta/m + \delta/(2m) \leq \delta$. That completes the induction.

PROOF OF THE THEOREM. Let Δ be the given measurable function of the observations. Let $L > 0$ be given and an integer $n \geq 1$ be given. Choose $a < b$ such that $b - a > 2nL$. Let F , having a bimodal density function with modes $\alpha < \beta$ satisfy Lemma 2.4 relative to the interval (a, b) and the number $\delta/(2n)$ where $\delta < 1$. Then $F(\alpha) < p < F(\beta)$. By modification of F on the interval (α, β) we may construct distribution functions G_1, \dots, G_{2n} such that G_i has the unique p -point γ_{p, G_i} which is the midpoint of the interval $(\alpha + (i - 1) \cdot (\beta - \alpha)/(2n), \alpha + i(\beta - \alpha)/(2n))$, $1 \leq i \leq 2n$. We may further suppose that

if $1 \leq i \leq 2n$ then G_i has a bimodal density function. As is easily shown,

$$(2.1) \quad \sum_{i=1}^{2n} P_F(|\Delta - \gamma_{p, G_i}| < L/2) \leq 1,$$

since the intervals involved are nonoverlapping. Further, by application of Lemma 2.4, if $1 \leq i \leq 2n$ then

$$(2.2) \quad |P_F(|\Delta - \gamma_{p, G_i}| < L/2) - P_{G_i}(|\Delta - \gamma_{p, G_i}| < L/2)| < \delta/(2n).$$

Addition of (2.1) and (2.2) gives

$$(2.3) \quad \sum_{i=1}^{2n} P_{G_i}(|\Delta - \gamma_{p, G_i}| < L/2) \leq \delta + 1 < 2.$$

From (2.3) we may conclude that for at least one index i , $P_{G_i}(|\Delta - \gamma_{p, G_i}| < L/2) < 1/n$. That completes the proof of the theorem.

REMARK. The proof of this paper may be carried out using density functions relative to a positive σ -finite measure μ defined on the Borel sets of the line provided $\mu(\{a\}) = 0$ for all real numbers a .

REFERENCES

- [1] BAHADUR, R. R. and SAVAGE, L. J. (1956). The nonexistence of certain statistical procedures in nonparametric problems. *Ann. Math. Statist.* **27** 1115-1122.
- [2] FARRELL, R. H. (1959). Sequentially determined bounded length confidence intervals. Ph.D. thesis, University of Illinois.
- [3] FARRELL, R. H. (1964). Asymptotic behavior of expected sample size in certain one sided tests. *Ann. Math. Statist.* **35** 36-72.
- [4] FARRELL, R. H. (1966). Bounded length confidence intervals for the p -point of a distribution function, III. *Ann. Math. Statist.* **37** 586-592.
- [5] WEISS, L. (1960). Confidence intervals of preassigned length for quantiles of unimodal populations. *Naval Res. Log. Quart.* **7** 251-256.