

AN EMPIRICAL DISTRIBUTION FUNCTION FOR SAMPLING WITH INCOMPLETE INFORMATION

BY MIRIAM AYER, H. D. BRUNK, G. M. EWING, W. T. REID,
AND EDWARD SILVERMAN

Sandia Corporation, University of Missouri, Northwestern University

Summary. For $i = 1, 2, \dots, n$, let N_i independent trials be made of an event with probability p_i , and suppose that the probabilities p_i are known to satisfy the inequalities $p_1 \geq p_2 \geq \dots \geq p_n$. Let a_i denote the number of successes in the i -th trial, and p_i^* the ratio a_i/N_i ($i = 1, 2, \dots, n$). Then the maximum likelihood estimates $\bar{p}_1, \dots, \bar{p}_n$ of the numbers p_1, \dots, p_n may be found in the following way. If $p_1^* \geq p_2^* \geq \dots \geq p_n^* \geq 0$, then $\bar{p}_i = p_i^*$, $i = 1, 2, \dots, n$. If $p_k^* \leq p_{k+1}^*$ for some k ($k = 1, 2, \dots, n - 1$), then $\bar{p}_k = \bar{p}_{k+1}$; the ratios $p_k^* = a_k/N_k$ and $p_{k+1}^* = a_{k+1}/N_{k+1}$ are then replaced in the sequence $p_1^*, p_2^*, \dots, p_n^*$ by the single ratio $(a_k + a_{k+1}) / (N_k + N_{k+1})$, obtaining an ordered set of only $n - 1$ ratios. This procedure is repeated until an ordered set of ratios is obtained which are monotone non-increasing. Then for each i , \bar{p}_i is equal to that one of the final set of ratios to which the original ratio a_i/N_i contributed. It is seen that this method of calculating the $\bar{p}_1, \dots, \bar{p}_n$ depends on a grouping of observations which might very well appeal to an investigator on purely intuitive grounds. It seems of interest to note that it yields the maximum likelihood estimates of the desired probabilities.

Particular examples of this situation are found in bio-assay [3] and in the proximity fuze problem discussed by M. Friedman ([1], Chapter 11).

The last section is devoted to a consistency property of the maximum likelihood estimators.

1. Introduction. In ordinary sampling one observes directly values of a random variable. There are, however, certain investigations, of which examples are to be found in a number of different fields, in which the result of each observation is not a sample value of the random variable being tested, but only a number, together with the information that the sample value is less than, or is greater than, that number. Bio-assay furnishes an example ([3]; for further references see [3], p. 416; [1], p. 352). Certain other examples occurring in the biological sciences have been suggested to the authors. Still another situation of this kind is mentioned by M. Friedman ([1], Chapter 11). Given a population of proximity fuzes, one is interested in the distribution of the random variable t , maximum distance from target at which a proximity fuze will operate. The result of a test of an individual proximity fuze is the distance of its nearest approach to the target and the information that it did or did not operate (we assume that the proximity fuze will not operate before reaching its point of closest approach to

Received March 8, 1954.

the target. We do not know whether or not this is in fact true of any actual proximity fuze; cf. [1], Chapter 11). The result of such an observation is therefore not a sample value of the random variable, \mathbf{t} , but rather a distance, t_0 , and the information that the sample value of \mathbf{t} corresponding to the particular proximity fuze is less than t_0 (if it did not operate) or at least t_0 (if it did operate).

Let $F(t) = \Pr\{\mathbf{t} < t\}$; $F(t)$ is the distribution function of the random variable \mathbf{t} . Let $p(t) = 1 - F(t) = \Pr\{\mathbf{t} \geq t\}$; $p(t)$ represents the probability that the fuze will operate if its minimum distance from the target is t . Suppose R fuzes are tested, and observed to pass within distances t_1, t_2, \dots, t_n of the target ($n \leq R$; several may have the same minimum distance from target); for convenience suppose the $\{t_i\}_1^n$ are arranged in increasing order. The R tests may be regarded as a set of R independent trials of events having probabilities $p_i = p(t_i)$ ($i = 1, 2, \dots, n$) of success (those observed at the same minimum distance from target having the same a priori probability of operating), if the term "success" is used to signify that the proximity fuze operated. The problem is to estimate the probabilities $\{p_i\}_1^n$ from the results of the R trials.

In a typical bio-assay situation, a large number of trials is made at each parameter value t_i ($i = 1, 2, \dots, n$). In such a situation the ratios, number of successes divided by number of trials, each determined for a particular parameter value, will with high probability be in monotone non-increasing order (assuming $t_1 \leq t_2 \leq \dots \leq t_n$). The "best" estimates of the probabilities are then these ratios, and if $\bar{p}(t)$ is a non-increasing function assuming these values at the points $\{t_i\}_1^n$ then $\bar{F}(t) = 1 - \bar{p}(t)$ is an obvious empirical distribution function. In other cases, such as that discussed above, one might expect a small number of trials corresponding to each parameter value, so that the average numbers of successes could not be expected to be in monotone order. It is for such situations that the maximum likelihood estimators of the probabilities $\{p(t_i)\}_1^n$ are determined in this paper. If $\{\bar{p}_i\}_1^n$ denotes the set of maximum likelihood estimates, and if $\bar{p}(t)$ is a monotone non-increasing function such that $\bar{p}(t_i) = \bar{p}_i$ ($i = 1, 2, \dots, n$) then $\bar{F}(t) = 1 - \bar{p}(t)$ will be termed an *empirical distribution function*.

In bio-assay situations it is often assumed that the random variable in question (perhaps after an elementary transformation) is normally distributed. Methods of probit analysis ([1], [2], [3]) have been developed for use with such an assumption. While it is true that an empirical distribution function may be useful in determining parameters of a normal distribution under such an assumption, the primary purpose of this paper is to present estimators of the probabilities $\{p(t_i)\}_1^n$ without reference to any assumption as to the distribution of the random variable being tested. These estimators are derived in section 2. The calculations required for their computation are extremely simple and rapid. In section 3, the consistency of the estimators is considered. A theorem is proved which states that the empirical distribution function, $\bar{F}(t) = 1 - \bar{p}(t)$, converges in probability to the distribution function $F(t)$ as the number of tests or trials becomes infinite in an appropriate way.

2. Maximum likelihood estimators of the probabilities. Let $a_i + b_i$ independent trials be made corresponding to the same parameter value, or observation point, t_i , of which a_i are successes ($i = 1, 2, \dots, n$). If $p_i = p(t_i)$ denotes the probability of success in a trial corresponding to the parameter value t_i ($i = 1, 2, \dots, n$), then the a priori probability of the event that, for each integer i , $1 \leq i \leq n$, a specified a_i of the trials will result in success is

$$(2.1) \quad \prod = \prod_{i=1}^n p_i^{a_i} (1 - p_i)^{b_i}.$$

Since $p(t)$ is non-increasing, the $\{p_i\}_1^n$ are known to satisfy the relations

$$(2.2) \quad 1 \geq p_1 \geq p_2 \geq \dots \geq p_n \geq 0.$$

The maximum likelihood estimates of $\{p_i\}_1^n$ are those numbers, $\{\bar{p}_i\}_1^n$, which maximize the probability \prod subject to the relations (2.2). (These estimates also maximize the probability,

$$\prod_{i=1}^n \binom{a_i + b_i}{a_i} p_i^{a_i} (1 - p_i)^{b_i},$$

that for each i , $1 \leq i \leq n$, there will be a_i successes among the $a_i + b_i$ trials at the observation point t_i).

In the context of the above discussion the numbers a_i and b_i ($i = 1, 2, \dots, n$) are non-negative integers. In section 3 they will be so regarded. However, the discussion of this section requires only that they be non-negative real numbers, such that $a_i + b_i > 0$ ($i = 1, 2, \dots, n$).

Let \mathfrak{P}_n denote the class of sets of real numbers $\{p_i\}_1^n$ satisfying the inequalities (2.2). The problem is to determine a set $\{\bar{p}_i\}_1^n$ in \mathfrak{P}_n affording a maximum value to \prod :

$$(2.3) \quad \prod_{i=1}^n \bar{p}_i^{a_i} (1 - \bar{p}_i)^{b_i} = \max_{\{p_i\} \in \mathfrak{P}_n} \prod_{i=1}^n p_i^{a_i} (1 - p_i)^{b_i}.$$

LEMMA 2.1. *There is a maximizing set $\{\bar{p}_i\}_1^n$.*

This follows immediately from the observation that the product is a continuous function of its arguments p_1, p_2, \dots, p_n and hence assumes its maximum on the closed, bounded set described by inequalities (2.2).

Set

$$(2.4) \quad p_i^* = a_i / (a_i + b_i) \quad (i = 1, 2, \dots, n)$$

THEOREM 2.1. *If $\{\bar{p}_i\}_1^n$ is a maximizing set, and if $\bar{p}_k > \bar{p}_{k+1}$ for some k , $1 \leq k \leq n$, then $p_k^* \geq \bar{p}_k > \bar{p}_{k+1} \geq p_{k+1}^*$. Also, $p_1^* \leq \bar{p}_1$, and $p_n^* \geq \bar{p}_n$.*

PROOF. We prove first that $p_k^* \geq \bar{p}_k$. The basis of the proof is the observation that the function $p^a(1 - p)^b$ increases for $0 \leq p < a/(a + b)$ and decreases for $a/(a + b) < p \leq 1$. Suppose $\bar{p}_k > p_k^*$. Choose $p'_k = \max(p_k^*, \bar{p}_{k+1})$. Then $\bar{p}_1 \geq \bar{p}_2 \geq \dots \geq \bar{p}_{k-1} > p'_k \geq \bar{p}_{k+1} \geq \dots \geq \bar{p}_n$, while

$$p_k'^{a_k} (1 - p'_k)^{b_k} > \bar{p}_k^{a_k} (1 - \bar{p}_k)^{b_k}.$$

This means that \prod is increased by replacing \bar{p}_k by p'_k (note that $\max \prod > 0$), contrary to (2.3). Therefore $p_k^* \geq \bar{p}_k$. Similarly $\bar{p}_{k+1} \geq p_{k+1}^*$. Hence $p_k^* \geq \bar{p}_k > \bar{p}_{k+1} \geq p_{k+1}^*$. The proof of the last statement of the theorem is similar.

For integers r, s , with $1 \leq r \leq s \leq n$, define

$$(2.5) \quad \begin{cases} \alpha(r, s) = \sum_{\nu=r}^s a_\nu, & \beta(r, s) = \sum_{\nu=r}^s b_\nu, \\ A(r, s) = \alpha(r, s)/[\alpha(r, s) + \beta(r, s)]. \end{cases}$$

THEOREM 2.2. For $1 \leq i \leq n$, we have

$$\begin{aligned} \bar{p}_i &= \min_{1 \leq r \leq i} \max_{i \leq s \leq n} A(r, s) = \max_{i \leq s \leq n} \min_{1 \leq r \leq i} A(r, s) \\ &= \min_{1 \leq r \leq i} \max_{r \leq s \leq n} A(r, s) = \max_{i \leq s \leq n} \min_{1 \leq r \leq s} A(r, s). \end{aligned}$$

The original proof, based on Theorem 2.1, is omitted. The reader is referred to the following paper for a simpler proof.

COROLLARY 2.1. The maximizing set $\{\bar{p}_i\}_1^n$ is unique. Each $\bar{p}_i (i = 1, 2, \dots, n)$ is determined uniquely by any of the formulas in Theorem 2.2.

Theorem 2.2 gives explicit formulas for the determination of the $\{\bar{p}_i\}$, but these are not recommended for calculation. Theorem 2.1 provides a means of calculating the maximizing set, $\{\bar{p}_i\}_1^n$, as outlined in the summary, which is very fast even for moderately large n .

The following interesting inequality was mentioned by a referee:

$$\sum_k (p_k^* - p_k)^2 (a_k + b_k) \geq \sum_k (\bar{p}_k - p_k)^2 (a_k + b_k).$$

Here p_k^* and \bar{p}_k are as defined above, while p_1, p_2, \dots, p_n is any set of numbers such that $1 \geq p_1 \geq p_2 \geq \dots \geq p_n \geq 0$. Indeed, one has

$$\sum_k (p_k^* - p_k)^2 (a_k + b_k) \geq \sum_k (\bar{p}_k - p_k)^2 (a_k + b_k) + \sum_k (p_k^* - \bar{p}_k)^2 (a_k + b_k),$$

as was shown by two of the authors, independently, in more general contexts, subsequent to the submission of the manuscript. These inequalities show that the numbers \bar{p}_k are, on the average (in an obvious sense), closer to the numbers p_k respectively than are the numbers p_k^* .

3. The consistency of the estimators. Let $F(t)$ be the distribution function of the random variable \mathbf{t} (see Section 1). The probability that \mathbf{t} will assume a value t or greater is given by $p(t) = 1 - F(t)$. The method discussed in Section 2 provides the maximum likelihood estimates, \bar{p}_i , of $p(t)$ at specified parameter values, or observation points, $t_i (i = 1, 2, \dots, n)$. Let $\bar{p}(t)$ denote any non-increasing function, $0 \leq \bar{p}(t) \leq 1$, assuming the values \bar{p}_i at the points $t_i (i = 1, 2, \dots, n)$, and $\bar{F}(t) = 1 - \bar{p}(t)$ an empirical distribution function associated with trials at the observation points $t_i (i = 1, 2, \dots, n)$. If the points t_1, \dots, t_n were to remain fixed and the number of trials at each to increase in-

definitely, it would follow from the strong law of large numbers that for $k = 1, 2, \dots, n$, p_k^* and \bar{p}_k converge with probability 1 to p_k . In the following theorem, however, neither n nor the points t_1, t_2, \dots, t_n , nor hence the probabilities p_k need remain fixed. For a fixed t_0 , the number of trials made at t_0 need not become infinite, nor need any at all be made at t_0 . We shall have $\bar{p}(t_0)$ near $p(t_0)$ with high probability if only enough trials are made at points near t_0 , even if only one trial is made at each point.

The following theorem of Kolmogorov (strong law of large numbers) will be useful in establishing such a result.

LEMMA 3.1. (Kolmogorov) *Let y_j be a sequence of independent random variables having expected values $E(y_j)$ and variances $V_j(j = 1, 2, \dots)$. Let ϵ be an arbitrary positive number, and M a positive integer. Then*

$$(3.1) \quad \Pr \left\{ \sup_{k \geq M} \left| \frac{1}{k} \sum_{j=1}^k [y_j - E(y_j)] \right| \leq \epsilon \right\} > 1 - \frac{16}{\epsilon^2} \left[\sum_{j=M}^{\infty} \frac{V_j}{j^2} + \frac{1}{4M^2} \sum_{j=1}^M V_j \right]$$

([4], p. 203).

THEOREM 3.1. *Let t_0 be a continuity point of the distribution function $F(t)$. Let ϵ, η be arbitrary positive numbers. Let t', t'' be chosen so that $t' < t_0 < t''$ and so that $|F(t) - F(t_0)| < \epsilon/2$ for $t' \leq t \leq t''$. Then*

$$(3.2) \quad \Pr\{|\bar{F}(t_0) - F(t_0)| < \epsilon\} > 1 - \eta$$

provided that at least N trials are made between t' and t_0 and at least N trials are made between t_0 and t'' , where N is chosen so that

$$(3.3) \quad \sum_{j=N}^{\infty} \frac{1}{j^2} + \frac{1}{4N} < \epsilon^2 \eta / 32.$$

PROOF. We shall prove first that $\Pr\{\bar{F}(t_0) > F(t_0) - \epsilon\} > 1 - \eta/2$ or

$$(3.4) \quad \Pr\{\bar{p}(t_0) < p(t_0) + \epsilon\} > 1 - \eta/2$$

provided that at least N trials are made between t' and t_0 . It can be shown similarly that $\Pr\{\bar{F}(t_0) < F(t_0) + \epsilon\} > 1 - \eta/2$, or

$$(3.5) \quad \Pr\{\bar{p}(t_0) < p(t_0) - \epsilon\} > 1 - \eta/2,$$

provided that at least N trials are made between t_0 and t'' . Inequality (3.2) follows from (3.4) and (3.5).

In order to establish (3.4), let $t^* = t_0$ if t_0 is an observation point. If not, let t^* denote the first observation point to the left of t_0 . Since $\bar{F}(t_0) \geq \bar{F}(t^*)$, or $\bar{p}(t_0) \leq \bar{p}(t^*)$, it suffices to prove

$$(3.6) \quad \Pr\{\bar{p}(t^*) < p(t_0) + \epsilon\} > 1 - \eta/2.$$

Let the observation points be $\{t_i\}(i = 1, 2, \dots, n)$ with $t_1 \leq t_2 \leq \dots \leq t_n$. Let t_m be the first observation point to the right of t' . Let M be the number of trials at observation points $t_m, t_{m+1}, \dots, t_u = t^*$. By hypothesis, $M \geq N$; that is, $\sum_{i=m}^u (a_i + b_i) \geq N$. Order the trials at observation points $t_m, t_{m+1}, \dots,$

t_n in the order of increasing t_i , ordering in an arbitrary way those occurring at the same observation point. Let $T_1, T_2, \dots, T_M, T_{M+1}, \dots, T_R$, where R is the total number of trials, denote the trials so ordered. Let $\{y_j\}$ denote the number of successes in the trial $T_j (j = 1, 2, \dots, R)$ so that $y_j = 1$ with probability $p(t_i)$ and $y_j = 0$ with probability $1 - p(t_i)$, where t_i is the observation point at which the trial T_j occurs. For $j > R$, let $\{y_j\}$ be independent random variables, each assuming the value $p(t_0)$ with probability 1. Set $s_k = \sum_{j=1}^k y_j$. By Theorem 2.2,

$$\bar{p}(t^*) = \bar{p}(t_u) = \min_{1 \leq r \leq u} \max_{u \leq s \leq n} A(r, s).$$

Hence

$$\bar{p}(t^*) \leq \max_{u \leq s \leq n} A(m, s)$$

(t_m is the first observation point to the right of t'). The symbol $A(m, s)$ represents the average number of successes in trials starting at t_m and terminating at t_s . Hence as s varies ($s \geq u$) these ratios form a subsequence of the sequence $s_k/k (k \geq M)$. This implies that

$$(3.7) \quad \bar{p}(t^*) \leq \sup_{k \geq M} s_k/k.$$

By Lemma 3.1,

$$\begin{aligned} \Pr \left\{ \sup_{k \geq M} \left[\frac{s_k}{k} - \frac{1}{k} \sum_{j=1}^k E(y_j) \right] \leq \frac{\epsilon}{2} \right\} &\geq \Pr \left\{ \sup_{k \geq M} \left| \frac{s_k}{k} - \frac{1}{k} \sum_{j=1}^k E(y_j) \right| \leq \frac{\epsilon}{2} \right\} \\ &> 1 - \frac{64}{\epsilon^2} \left[\sum_{j=M}^{\infty} V_j/j^2 + \frac{1}{4M^2} \sum_{j=1}^M V_j \right]. \end{aligned}$$

But $V_j = \text{Var}(y_j) = p(t_i)[1 - p(t_i)] \leq \frac{1}{4}$, t_i being the observation point at which the trial T_j occurs. Hence by hypothesis (3.3),

$$\Pr \left\{ \sup_{k \geq M} \left[\frac{s_k}{k} - \frac{1}{k} \sum_{j=1}^k E(y_j) \right] \leq \frac{\epsilon}{2} \right\} > 1 - \frac{16}{\epsilon^2} \left[\sum_{j=M}^{\infty} \frac{1}{j^2} + \frac{1}{4M} \right] > 1 - \eta/2,$$

since $M \geq N$. Further, if $1 \leq j \leq R$, then $E(y_j) = p(t_i) < p(t_0) + \epsilon/2$; if $j > R$, then $E(y_j) = p(t_0)$. Hence

$$\Pr \left\{ \sup_{k \geq M} s_k/k < p(t_0) + \epsilon \right\} > 1 - \eta/2.$$

By (3.7) it then follows that

$$(3.6) \quad \Pr \{ \bar{p}(t^*) < p(t_0) + \epsilon \} > 1 - \eta/2.$$

The proof of Theorem 3.1 is completed as indicated immediately following its statement.

REFERENCES

- [1] E. EISENHART, M. W. HASTAY, W. A. WALLIS (Editors), *Techniques of Statistical Analysis*, McGraw-Hill, New York, 1947.
- [2] D. J. FINNEY, *Probit Analysis*, Cambridge University Press, 2nd edition, 1951.
- [3] C. H. GOULDEN, *Methods of Statistical Analysis*, John Wiley and Sons, 2nd edition, New York, 1952.
- [4] J. V. USPENSKY, *Introduction to Mathematical Probability*, McGraw-Hill, New York, 1937.