

A SEQUENTIAL DECISION PROCEDURE FOR CHOOSING ONE OF k HYPOTHESES CONCERNING THE UNKNOWN MEAN OF A NORMAL DISTRIBUTION¹

BY EDWARD PAULSON

Queens College

1. Summary. A sequential procedure is given for deciding to which of k non-overlapping intervals the unknown mean θ belongs which satisfies the requirement that the probability of making an incorrect decision is less than some pre-assigned value α . The sequential procedure is worked out explicitly for the following two cases: (1) when θ is the mean of a normal distribution with a known variance, and (2) when θ is the mean of a normal distribution with an unknown variance. A brief discussion is also given of a related but apparently new problem, to find a sequential procedure which will simultaneously select one of the k intervals and also yield a confidence interval for θ of a specified width.

2. Introduction. Let X denote a random variable with probability density function

$$f(x) = [\sigma(2\pi)^{\frac{1}{2}}]^{-1} \exp [-(x - \theta)^2/2\sigma^2].$$

Let $\{I_j\}$ ($j = 1, 2, \dots, k$) denote k non-overlapping intervals whose union $\cup_{j=1}^k I_j$ is the real line. A problem that seems to be of considerable practical interest is that of deciding on the basis of a sample $\{X_i\}$ of independent measurements on X to which of the k intervals the unknown mean θ belongs. For example, a manufacturer might want to make k different decisions concerning the selling price of his output, according to which of the k intervals contains the unknown mean. Another application concerns classification problems which only involve a single characteristic, where an anthropologist might want to decide which of k known populations

$$f_j(x) = [\sigma_j(2\pi)^{\frac{1}{2}}]^{-1} \exp [-(x - \theta_j)^2/2\sigma_j^2]$$

($j = 1, 2, \dots, k$) is "closest" to

$$f(x) = [\sigma(2\pi)^{\frac{1}{2}}]^{-1} \exp [-(x - \theta)^2/2\sigma^2].$$

Assuming for simplicity that $\theta_1 < \theta_2 < \dots < \theta_k$, then guided by practical or theoretical considerations we can introduce $k - 1$ numbers $\lambda_1, \lambda_2, \dots, \lambda_{k-1}$ (ordinarily $\lambda_j = (\theta_j + \theta_{j+1})/2$ if $\sigma_1 = \sigma_2 = \dots = \sigma_k = \sigma$) so that the problem is reduced to deciding which of the k intervals

$$(-\infty, \lambda_1], (\lambda_1, \lambda_2], (\lambda_2, \lambda_3], \dots, (\lambda_{k-1}, \infty)$$

contains θ .

Received July 26, 1962.

¹ This research was supported by the National Science Foundation under Grant NSF-G23665.

A sequential solution to the problem specified above was given by Sobel and Wald [1] for the special case when $k = 3$ and σ was assumed known. A sequential solution for the general case is given in the present paper using a different approach. No claim is made that the solution given here has any optimum properties, but nevertheless it is hoped that the results may be useful in practice.

3. The case when σ is known. Let $\{X_i\}$ $i = 1, 2, \dots$ be a series of independent and identically distributed random variables with common probability density

$$f(x) = [\sigma(2\pi)^{\frac{1}{2}}]^{-1} \exp [-(x - \theta)^2/2\sigma^2],$$

where σ is assumed known a priori. Let $\lambda_1 < \lambda_2 < \lambda_3 \dots < \lambda_{k-1}$, let I_1 denote the interval $(-\infty, \lambda_1]$, for $s = 2, 3, \dots, k-1$ let I_s denote the half open interval $(\lambda_{s-1}, \lambda_s]$, and let I_k denote the interval (λ_{k-1}, ∞) . Let D_j ($j = 1, 2, \dots, k$) denote the decision that θ is an element of I_j .

Following the formulation of the problem given in [1], we assume we can specify on the basis of practical considerations an indifference zone (c_j, d_j) about each point λ_j , in which it is a matter of indifference whether decisions D_j or D_{j+1} are made. Let $W(\theta, D_j)$ be a function taking values 0 or 1 which represents the error made in selecting D_j when θ is the true parameter. We assume the function $W(\theta, D_j)$ is given as follows:

$$\begin{aligned} W(\theta, D_1) &= 0 && \text{if } -\infty < \theta < d_1 \\ &= 1 && \text{otherwise.} \end{aligned}$$

For $s = 2, 3, \dots, k-1$

$$\begin{aligned} W(\theta, D_s) &= 0 && \text{if } c_{s-1} < \theta < d_s \\ &= 1 && \text{otherwise.} \end{aligned}$$

$$\begin{aligned} W(\theta, D_k) &= 0 && \text{if } c_{k-1} < \theta < \infty \\ &= 1 && \text{otherwise.} \end{aligned}$$

We now proceed to find a sequential procedure for choosing one of the k decisions (D_1, D_2, \dots, D_k) so that the probability of making an error does not exceed a preassigned value α for all values of θ . We will repeatedly make use of the following result (see for example Section 2.1 of [2]: if Y_1, Y_2, \dots is a sequence of independent observations on a random variable Y with $E(Y) < 0$, then if $a > 0$ and $r = 1, 2, 3, \dots$

$$(1) \quad P \left[\sum_{i=1}^r y_i > a \text{ for at least one value of } r \right] \leq e^{-t_0 a},$$

where $t_0 \neq 0$ satisfies the relation $E[e^{t_0 Y}] = 1$. Applying this result to $Y_i = (X_i - \theta - d)/\sigma^2$ where $d > 0$, we find that

$$(2) \quad P \left[\sum_{i=1}^r [(X_i - \theta - d)/\sigma^2] > a \text{ for at least one value of } r \right] \leq e^{-2da},$$

which is equivalent to

$$(3) \quad P[\bar{x}_r - d - a\sigma^2/r > \theta \text{ for at least one value of } r] \leq e^{-2da}$$

where $\bar{x}_r = \sum_{i=1}^r X_i/r$. Applying (1) again with Y_i now equal to

$$(-X_i + \theta - d)/\sigma^2$$

we find in a similar manner that

$$(4) \quad P[\bar{x}_r + d + a\sigma^2/r < \theta \text{ for at least one value of } r] \leq e^{-2da}.$$

Now setting $e^{-2da} = \alpha/2$ and applying the inequality

$$P[A \cap B] \geq 1 - P(\bar{A}) - P(\bar{B})$$

to (3) and (4), we find that

$$(5) \quad P[\bar{x}_r - d - \sigma^2 \log(2/\alpha)/2dr \leq \theta \leq \bar{x}_r + d + \sigma^2 \log(2/\alpha)/2dr \text{ for every } r, r = 1, 2, \dots] \geq 1 - \alpha$$

where all logarithms are to the base e .

For the application of (5) it is desirable to choose d so as to minimize the value of r required to have

$$2\{d + \sigma^2 \log(2/\alpha)/2dr\} = \Delta.$$

A routine calculation shows that this is accomplished when $d = \Delta/4$.

The required sequential procedure can be obtained on the basis of the result in (5) as follows. First, let $\Delta_j = d_j - c_j$, the length of the j th indifference zone, let $\Delta = \min(\Delta_1, \Delta_2, \dots, \Delta_{k-1})$ and let

$$u_r = \max_i (1 \leq i \leq r) \{ \bar{x}_i - \Delta/4 - 2\sigma^2 \log(2/\alpha)/i\Delta \},$$

$$v_r = \min_i (1 \leq i \leq r) \{ \bar{x}_i + \Delta/4 + 2\sigma^2 \log(2/\alpha)/i\Delta \}.$$

We now describe the sequential procedure. At the r th stage of the experiment ($r = 1, 2, \dots$) we obtain X_r and then calculate u_r and v_r on the basis of the observed measurements X_1, X_2, \dots, X_r . The experiment is terminated at the r th stage if either $u_r > v_r$ or if $u_r \leq v_r$ and the interval (u_r, v_r) falls entirely within one of the k intervals $I_1' = (-\infty, d_1), I_2' = (c_1, d_2), I_3' = (c_2, d_3), \dots, I_{k-1}' = (c_{k-2}, d_{k-1}), I_k' = (c_{k-1}, \infty)$. If the experiment terminates at the r th stage with $u_r > v_r$ (the probability of this will ordinarily be quite small), we make the decision corresponding to the interval among $\{I_j\}$ to which \bar{x}_r belongs. If the experiment terminates at the r th stage with $u_r \leq v_r$, we make a decision corresponding to the interval among $\{I_j'\}$ in which (u_r, v_r) lies, so that if (u_r, v_r) say falls in I_i' , we select D_i , while if (u_r, v_r) falls inside two intervals, say I_i' and I_{i+1}' , we can choose between D_i and D_{i+1}' by throwing a coin. If no decision is

reached at the r th stage of the experiment, we go on to the $(r + 1)$ st stage, obtain X_{r+1} , and repeat the above procedure.

Let n denote the stage where the experiment terminates. It is clear that the experiment must have terminated on or before the j th stage if

$$2[\Delta/4 + 2\sigma^2 \log(2/\alpha)/j\Delta] < \Delta,$$

so that we must have $n < 1 + 8\sigma^2 \log(2/\alpha)/\Delta^2$. The sequential procedure given is therefore closed, which is a very useful feature in most applications.

From the definition of the error function $W(\theta, D_j)$ and the specification of the sequential procedure, it is clear that no error can be committed when θ is really contained in (u_n, v_n) . Since it follows from (5) that $P[u_n \leq \theta \leq v_n] \geq 1 - \alpha$ we can therefore conclude that the probability the sequential procedure will lead to an error is $\leq \alpha$ for all θ .

Although a detailed investigation has not been made, it seems very likely that when $k = 3$ the sequential procedure given here is less efficient than the Sobel-Wald procedure discussed in [1]. However, assuming this to be the case, there is some compensation because (1) the sequential procedure given here is closed and (2) when the sequential procedure given here terminates, some extra information is available, since in addition to making a decision regarding the interval in which θ lies, we can also assert that $u_n \leq \theta \leq v_n$ with probability $\geq 1 - \alpha$.

The last remark in the above paragraph suggests a modification of the original problem which might be of interest in itself. The modified problem is as follows: to find a sequential procedure which will simultaneously select one of the decisions (D_1, D_2, \dots, D_k) with the probability of an error $\leq \alpha$, and also provide a confidence interval for θ of width $\leq W$ with confidence coefficient $\geq 1 - \alpha$. (W is a constant which is assumed to be $> \Delta$).

To solve this modified problem, we use the sequential procedure described before with the modification that we continue sampling until not only (u_r, v_r) falls inside one of the intervals $\{I'_j\}$ but in addition $v_r - u_r \leq W$. (The minor complication caused by the fact that it is possible to have $u_r > v_r$ can be handled by taking $\bar{x}_r \pm W/2$ as the confidence limits for θ when this unlikely event occurs.) It is clear that with this modification the probability is still $\leq \alpha$ of making an incorrect decision and the probability is $\geq 1 - \alpha$ that $u_n \leq \theta \leq v_n$ and that $0 \leq v_n - u_n \leq W$.

4. The case when σ is unknown. We now start the experiment by first taking a preliminary sample of n_0 measurements X_1, X_2, \dots, X_{n_0} . Let $f = n_0 - 1$ and let

$$s^2 = \sum_{i=1}^{n_0} (X_i - \bar{x}_{n_0})^2 / f.$$

Let $\{Z_i\}$ be a sequence of independent and normally distributed random variables with mean $-d/\sigma^2$ and variance $1/\sigma^2$ which is independent of the sequence $\{X_i\}$.

Making use of the fact that \bar{x}_{n_0} is independent of s^2 , we obtain

$$\begin{aligned}
 & P \left[\sum_{i=1}^r [(X_i - \theta - d)/s^2] > a \text{ for at least one } r, \quad r = n_0, n_0 + 1, \dots \mid s^2 \right] \\
 &= P \left[\sum_{i=1}^r [(X_i - \theta - d)]/\sigma^2 > a (s^2/\sigma^2) \text{ for at least one } r, \right. \\
 (6) \quad & \left. r = n_0, n_0 + 1, \dots \mid s^2 \right] \\
 &= P \left[\sum_{i=1}^r Z_i \geq a(s^2/\sigma^2) \text{ for at least one } r, \quad r = n_0, n_0 + 1, \dots \mid s^2 \right] \\
 &< P \left[\sum_{i=1}^r Z_i \geq a(s^2/\sigma^2) \text{ for at least one } r, \quad r = 1, 2, \dots \mid s^2 \right] \leq e^{-2ads^2/\sigma^2}.
 \end{aligned}$$

Therefore

$$\begin{aligned}
 (7) \quad & P \left[\sum_{i=1}^r [(X_i - \theta - d)/s^2] > a \text{ for at least one } r, \quad r = n_0, n_0 + 1, \dots \right] \\
 & < E \{ e^{-2ads^2/\sigma^2} \} = [1 + 4ad/f]^{-f/2}.
 \end{aligned}$$

If we set $d = \Delta/4$ and let $\bar{a} = (f/\Delta)[(2/\alpha)^{2/f} - 1]$ denote the value of a for which $[1 + 4ad/f]^{-f/2} = \alpha/2$, we obtain from (7)

$$(8) \quad P[\bar{x}_r - \Delta/4 - \bar{a}s^2/r > \theta \text{ for at least one } r, \quad r = n_0, n_0 + 1, \dots] < \alpha/2.$$

Starting with $(-X_i + \theta - d)/s^2$, we find in a similar manner that

$$(9) \quad P[\bar{x}_r + \Delta/4 + \bar{a}s^2/r < \theta \text{ for at least one } r, \quad r = n_0, n_0 + 1, \dots] < \alpha/2.$$

Therefore combining (8) and (9),

$$\begin{aligned}
 (10) \quad & P[\bar{x}_r - \Delta/4 - \bar{a}s^2/r \leq \theta \leq \bar{x}_r + \Delta/4 + \bar{a}s^2/r \\
 & \text{for all values of } r, \quad r = n_0, n_0 + 1, \dots] > 1 - \alpha.
 \end{aligned}$$

Let

$$u'_r = \max_j (n_0 \leq j \leq r) \{ \bar{x}_j - \Delta/4 - \bar{a}s^2/j \},$$

and

$$v'_r = \min_j (n_0 \leq j \leq r) \{ \bar{x}_j + \Delta/4 + \bar{a}s^2/j \}.$$

Now we can specify the sequential procedure when σ is unknown. We start by taking n_0 observations, compute s^2 , and for $r = n_0, n_0 + 1, \dots$ we then use the sequential procedure described in Section 3 with u_r and v_r replaced by u'_r and v'_r . It follows from (10) that the probability of making an error is still $< \alpha$ for all θ .

The practical value of the procedure in Section 4 is somewhat limited by the

lack of an efficient rule for selecting n_0 . This appears to be a fairly important unsolved problem.

REFERENCES

- [1] SOBEL, M. and WALD, A. (1949). A sequential decision procedure for choosing one of three hypotheses concerning the unknown mean of a normal distribution. *Ann. Math. Statist.* **20** 502–22.
- [2] BARTLETT, M. S. (1955). *An Introduction to Stochastic Processes*. Cambridge Univ. Press.