

## AN APPROXIMATION TO THE SAMPLE SIZE IN SELECTION PROBLEMS<sup>1</sup>

BY EDWARD J. DUDEWICZ

*The University of Rochester*

**0. Summary.** Let  $f(\mathbf{x} | P_1)$  be the pdf of a  $(k - 1)$ -dimensional normal distribution with zero means, unit variances, and correlation matrix  $P_1$ . Consider the integral, for  $\delta > 0$ ,

$$(1) \quad \int_{-\delta}^{\infty} \cdots \int_{-\delta}^{\infty} f(\mathbf{x} | P_1) dx \cdots dx_{k-1} = \alpha(\delta), \quad \text{say.}$$

Assume that no element of  $P_1$  is a function of  $\delta$ . Note that  $\alpha(\delta)$  is an increasing function of  $\delta$  and  $\alpha(\delta) \rightarrow 1$  as  $\delta \rightarrow \infty$ . The problem is to obtain an approximation to  $\delta$ , for a large specified value,  $\alpha$ , of  $\alpha(\delta)$ . This is given by the theorem of Section 1.

This result is used to obtain approximations to the sample size in a selection procedure of Bechhofer and in a problem of selection from a multivariate normal population. The closeness of the approximation is illustrated for the procedure of Bechhofer (Table 1).

### 1. The Approximation to $\delta$ .

**THEOREM.** For large  $\alpha$  (near 1), an approximation to  $\delta$ , which satisfies the equation

$$(2) \quad \int_{-\delta}^{\infty} \cdots \int_{-\delta}^{\infty} f(\mathbf{x} | P_1) dx_1 \cdots dx_{k-1} = \alpha,$$

is  $\delta^2 \sim -2 \log_e (1 - \alpha)$ . The ratio tends to 1 as  $\alpha \rightarrow 1$ . This approximation is independent of  $k$ .

**PROOF.**

$$(3) \quad \begin{aligned} \alpha(\delta) &= \int_{-\delta}^{\infty} \cdots \int_{-\delta}^{\infty} f(\mathbf{x} | P_1) dx_1 \cdots dx_{k-1} \\ &= P[\bigcap_{i=1}^{k-1} \{Z_i > -\delta\}] = P[\bigcap_i E_i], \quad \text{say.} \end{aligned}$$

Then,

$$(4) \quad 1 - \alpha(\delta) = P[\mathbf{U}_i \bar{E}_i].$$

Expressing the union as a disjoint union, we obtain

$$(5) \quad 1 - \alpha(\delta) = \sum_{i=1}^{k-1} P[E_1 E_2 \cdots \bar{E}_i].$$

Now

$$P[E_1 E_2 \cdots \bar{E}_i] = \int_{-\delta}^{\infty} \cdots \int_{-\delta}^{\infty} \int_{-\delta}^{-\delta} f_i(\mathbf{x} | P_2) dx_1 \cdots dx_{i-1} dx_i,$$

Received 2 February 1968; revised 13 September 1968.

<sup>1</sup> The author wishes to acknowledge the support of an NSF Summer Fellowship and of a Research Assistantship supported by the ONR on Contract Nonr 401 (53) at Cornell University, and of ONR Contract N00014-68A-0091 at The University of Rochester. "This research is supported in whole or in part by The Center for Naval Analyses of the University of Rochester. Such support does not imply endorsement of the content by The Navy." Revision of this paper was supported by NSF Grant No. GP-8958.

where  $f_i(\mathbf{x} | P_2)$  is the pdf of an  $i$ -dimensional normal distribution with zero means, unit variances, and correlation matrix obtained from  $P_1$  by deleting the last  $k - i$  rows and columns. The  $j$ th of the  $i - 1$  integrals on  $(-\delta, \infty)$  can be replaced by some multiplier  $C_j(i)$  between 0 and  $(\pi^{k-1}k)^{\frac{1}{2}}$  which tends to a limit as  $\delta \rightarrow \infty$  ( $j = 1, \dots, i - 1$ ), and using problem 1 of Feller (1957), p. 179, the remaining integral on  $(-\infty, -\delta)$ , which equals the integral on  $(\delta, \infty)$ , can be approximated as

$$1 - \Phi(\delta) \sim \delta^{-1}(2\pi)^{-\frac{1}{2}}e^{-\delta^2/2}.$$

Letting  $C_i = C_1(i) \cdots C_{i-1}(i)$ ,

$$P[E_1E_2 \cdots \bar{E}_i] \sim C_i(\delta(2\pi)^{\frac{1}{2}})^{-1}e^{-\delta^2/2},$$

and letting  $\bar{C} = (C_1 + \cdots + C_{k-1})(k - 1)^{-1}$ ,  $\bar{C} = \bar{C}(\delta) \rightarrow$  a limit as  $\delta \rightarrow \infty$  and an approximation to  $\delta$ , such that  $\alpha(\delta) = \alpha$ , can be obtained from the equation

$$1 - \alpha \sim (k - 1)\bar{C}(\delta(2\pi)^{\frac{1}{2}})^{-1}e^{-\delta^2/2}.$$

Hence  $\delta^2 \sim -2 \log_e (1 - \alpha)$ .

**2. The approximation to sample size.** Ranking and selection procedures have been receiving increasing interest recently, one reason being that they furnish a method of sample size selection which is appropriate in many situations. However, the tables available for this use (e.g. those of Bechhofer (1954) and Teichrow (1955)) are not extensive enough to cover all situations which arise in practice, e.g. in drug screening where the number of populations may be very large. Also, these tables do not apply if certain independence assumptions are violated. The theorem of Section 1 provides an approximation to the sample size in a selection procedure of Bechhofer and in a multivariate situation.

Consider the following *problem*: Given  $k$  populations  $\pi_1, \dots, \pi_k$  the observations from which are known to be identically distributed except for a location parameter (i.e., an observation  $X_i$  from population  $\pi_i$  has cdf  $F(x - \nu_i)$ , where  $F(\cdot)$  may be either known or unknown), select any one of the (at least one) populations associated with  $\nu_{[k]} = \max(\nu_1, \dots, \nu_k)$ . For the case in which  $F(\cdot)$  is known to be normal with a known variance  $\sigma^2$ , Bechhofer (1954) has suggested use of the *means procedure*: Take  $N$  independent vectors  $\underline{X}_j = (X_{1j}, \dots, X_{kj})$ ,  $j = 1, \dots, N$ , where  $X_{ij}$  denotes the  $j$ th observation from the  $i$ th population  $\pi_i$ ; choose the population associated with the largest  $\sum_{j=1}^N X_{ij}$  ( $i = 1, \dots, k$ ) as being associated with  $\nu_{[k]}$ . Let  $\{\lambda^*, P^*\}$  ( $0 < \lambda^* < \infty, 1/k < P^* < 1$ ) be two specified constants, and denote the ranked means by  $\nu_{[1]} \leq \dots \leq \nu_{[k]}$ . Then  $N$  is to be set as the smallest sample size which guarantees the following *probability requirement*: We are to select the population associated with  $\nu_{[k]}$ , i.e. we are to make a *correct selection* (CS), with probability  $P(\text{CS}) \geq P^*$  whenever  $\nu_{[k]} - \nu_{[k-1]} \geq \lambda^* \sigma$ .

It is known (see Bechhofer (1954), p. 23) that, if the means procedure is used for this problem in the normal case, then subject to  $\nu_{[k]} - \nu_{[k-1]} \geq \lambda^* \sigma$  the  $P(\text{CS})$

is minimized over  $\nu_{[1]}, \dots, \nu_{[k]}$  when

$$(6) \quad \nu_{[1]} = \dots = \nu_{[k-1]}, \nu_{[k]} - \nu_{[k-1]} = \lambda^* \sigma,$$

the *least favorable configuration* (LFC) of the population parameters; call this minimal  $P(\text{CS})$  by  $\alpha$ .

**THEOREM.** As  $\alpha \rightarrow 1$ ,

$$(7) \quad N \sim -4(\lambda^*)^{-2} \log_e (1 - \alpha),$$

the ratio tending to 1 as  $N \rightarrow \infty$  due to having  $\alpha \rightarrow 1$ .

**PROOF.** From Bechhofer (1954), p. 20, eq. (13), we know that

$$(8) \quad \alpha = \int_{-\lambda^*(\frac{1}{2}N)^{\frac{1}{2}}}^{\infty} \dots \int_{-\lambda^*(\frac{1}{2}N)^{\frac{1}{2}}}^{\infty} f(\mathbf{x} | P_1) dx_1 \dots dx_{k-1},$$

where  $P_1$  is the  $(k - 1) \times (k - 1)$  correlation matrix with  $\rho_{ij} = \frac{1}{2}(i \neq j)$ . The theorem follows from (2).

**3. Numerical study of the approximation.** We can compare the  $N_A$  calculated via (7) with  $N_T$  from the tables of Bechhofer (1954) for moderate  $k$ , obtaining Table 1. (We use  $\lambda^* = 1$  for convenience.) The tables of Gupta (1963), p. 810, and of Teichroew (1955) are useful in making comparisons for large  $k$  and for high  $P^*$ , respectively.

Note that the approximation (7) appears to overestimate  $N$ , thus being somewhat conservative, for  $k \leq 10$ . If we had  $\lambda^* \neq 1$ , each entry ( $N_T$  as well as  $N_A$ ) would be divided by  $(\lambda^*)^2$ . Thus, the ratio (which by derivation approaches 1 as  $N \rightarrow \infty$  due to having  $\alpha \rightarrow 1$ ) is independent of  $\lambda^*$ . Although it does not appear to be the case that  $|N_A - N_T| \rightarrow 0$  (indeed the difference could conceivably be unbounded) the approach of the ratio to 1 appears fairly rapid.

**4. Extension to multivariate normal selection.** Use of the special nature of the correlation matrix  $P_1$  in (8) has not been made. Thus, result (7) is valid whenever a reduction of  $\alpha$  to form (8) is possible, provided only that  $P_1$  does not depend on  $N$ .

Now, consider the following *problem*: Given one  $k$ -variate population  $\Pi$  the observations from which are  $k$ -variate normal, select any one of the (at least one) factors associated with  $\nu_{[k]}$ . Suppose that the *means procedure* stated in Section 2

TABLE 1

$\alpha$		0.90	0.95	0.97	0.99	0.995	0.999
	$N_A$	9.2104	11.9829	14.0262	18.4207	21.1933	27.6310
$k = 3$	$N_T$	4.9738	7.3446	9.1397	13.0849	15.6159	21.5760
	$N_T/N_A$	.54	.61	.65	.71	.74	.78
$k = 5$	$N_T$	6.7584	9.3342	11.2419	15.3633	17.9725	24.0580
	$N_T/N_A$	.73	.78	.80	.83	.85	.87
$k = 10$	$N_T$	8.8977	11.6841	13.7122	18.0251	20.7234	26.9537
	$N_T/N_A$	.97	.98	.98	.98	.98	.98

is used. We then observe  $N$  independent  $k$ -variate normal vectors  $\underline{X}_j = (X_{1j}, \dots, X_{kj}), j = 1, \dots, N$ , and select the factor associated with the largest  $\bar{X}_i = \sum_{j=1}^N X_{ij}/N (i = 1, \dots, N)$  as being associated with  $\nu_{[k]}$ .

Denote the  $k$ -variate normal distribution of any  $\underline{X}_j (j = 1, \dots, N)$  by  $N(\mathbf{v}_1, \Sigma_1)$  where  $\mathbf{v}_1 = (\nu_1, \dots, \nu_k)$  is unknown and  $\Sigma_1$  is some covariance matrix, say  $\Sigma_1 = (\sigma_{ij})$  for  $i, j = 1, \dots, k$ . Let round brackets about a subscript denote the quantity associated with that one of the ranked means; thus,  $\underline{X}_{(j)} (j = 1, \dots, N)$  is  $N(\mathbf{v}_{(1)}, \Sigma_{(1)})$  where  $\mathbf{v}_{(1)} = (\nu_{(1)}, \dots, \nu_{(k)})$  and  $\Sigma_{(1)} = (\sigma_{(i)(j)})$  for  $i, j = 1, \dots, k$ . Now, our development will be a generalization of Bechhofer (1954), p. 20. Then (assuming  $\nu_{[k]} - \nu_{[k-1]} > 0$ )

$$\begin{aligned}
 P(\text{CS}) &= P[\bar{X}_{(1)} < \bar{X}_{(k)}, \dots, \bar{X}_{(k-1)} < \bar{X}_{(k)}] \\
 (9) \qquad &= P[\bar{X}_{(k)} - \bar{X}_{(1)} > 0, \dots, \bar{X}_{(k)} - \bar{X}_{(k-1)} > 0] \\
 &= P[Y_1 > 0, \dots, Y_{k-1} > 0], \quad \text{say,}
 \end{aligned}$$

where

$$(\bar{X}_{(1)}, \dots, \bar{X}_{(k)}) \text{ is } N(\mathbf{v}_{(1)}, N^{-1}\Sigma_{(1)}).$$

Now,  $(Y_1, \dots, Y_{k-1})$  is  $(k - 1)$ -variate normal say  $N(\delta, \Sigma_2)$  with  $\delta_i = EY_i = \nu_{[k]} - \nu_{[i]}$  and

$$\begin{aligned}
 (10) \qquad \sigma^2(Y_i) &= N^{-1}(\sigma_{(k)}^2 + \sigma_{(i)}^2 - 2\sigma_{(i)(k)}) \\
 \sigma(Y_i Y_j) &= N^{-1}(\sigma_{(k)}^2 - \sigma_{(i)(k)} - \sigma_{(j)(k)} + \sigma_{(i)(j)}).
 \end{aligned}$$

Now, (9) is minimized subject to  $\nu_{[k]} - \nu_{[k-1]} \geq \lambda^* \sigma$  (now  $\sigma > 0$  is simply some constant) by the LFC (6). Thus,

$$\begin{aligned}
 (11) \qquad \alpha &= P[(Y_i - \lambda^* \sigma)/\sigma(Y_i) > -\lambda^* \sigma/\sigma(Y_i), \quad i = 1, \dots, k - 1] \\
 &= P[Z_i > -\lambda^* \sigma/\sigma(Y_i), \quad i = 1, \dots, k - 1], \text{ say.}
 \end{aligned}$$

$(Z_1, \dots, Z_{k-1})$  is  $(k - 1)$ -variate normal  $N(0, \Sigma_3)$  with

$$(12) \qquad \sigma^2(Z_i) = 1, \quad \sigma(Z_i Z_j) = \sigma(Y_i Y_j)/\sigma(Y_i)\sigma(Y_j).$$

Therefore,

$$\begin{aligned}
 (13) \qquad \alpha &= |\Sigma_3|^{-\frac{1}{2}}/(2\pi)^{\frac{1}{2}(k-1)} \int_{(-\lambda^* \sigma/\sigma(Y_1))}^{\infty} \dots \int_{(-\lambda^* \sigma/\sigma(Y_{k-1}))}^{\infty} \exp(-\frac{1}{2}x' \Sigma_3^{-1}x) dx_1 \dots dx_{k-1} \\
 &\geq |\Sigma_3|^{-\frac{1}{2}}/(2\pi)^{\frac{1}{2}(k-1)} \int_{(-\lambda^* \sigma N^{\frac{1}{2}}/2^{\frac{1}{2}} \sigma_M)}^{\infty} \dots \int_{(-\lambda^* \sigma N^{\frac{1}{2}}/2^{\frac{1}{2}} \sigma_M)}^{\infty} \exp(-\frac{1}{2}x' \Sigma_3^{-1}x) dx_1 \dots dx_{k-1},
 \end{aligned}$$

where

$$(14) \qquad \sigma_M = \max_{1 \leq i \leq k-1} [\frac{1}{2}(\sigma_{(k)}^2 + \sigma_{(i)}^2) - \sigma_{(i)(k)}]^{\frac{1}{2}}.$$

Since  $\Sigma_3$  is independent of  $N$ , our proof of (7) shows that the minimal  $N$  needed to make the lower bound in (13) equal to  $\alpha$ , say  $\bar{N}$ , is such that we have the

THEOREM. As  $\alpha \rightarrow 1$ ,

$$(15) \quad \bar{N} \sim -4(\lambda^*)^{-2}(\sigma/\sigma_M)^{-2} \log_e (1 - \alpha),$$

the ratio tending to 1 as  $\bar{N} \rightarrow \infty$  due to having  $\alpha \rightarrow 1$ .

Due to its construction,  $\bar{N}$  overprotects against an incorrect selection. Even if one has *no knowledge whatever* of  $\sigma_M$ , result (15) allows one to obtain an approximation to  $\bar{N}$ , if one sets  $\sigma = \sigma_M$ , i.e. if one makes his probability requirement in terms of multiples of the (unknown)  $\sigma_M$ .

If one has *special knowledge* then  $\sigma_M$  may sometimes be evaluated. E.g., if  $\sigma_1^2 = \dots = \sigma_k^2 = \sigma_0^2$  (say) and  $\sigma_{(i)(k)} \geq -B$  ( $B \geq 0$ ) for  $i = 1, \dots, k$  with at least one equality, then  $\sigma_M = (\sigma_0^2 + B)^{\frac{1}{2}}$ . In this case the "best" such  $B$  (in terms of  $\bar{N}$ ) is  $B = 0$ . (If one knew only  $\sigma_{(i)(k)} \geq C$  for  $i = 1, \dots, k$ , then one would obtain  $\sigma_M \geq (\sigma_0^2 - C)^{\frac{1}{2}}$ .)

The upper bound  $\bar{N}$  will be "exact" as  $N$  is in (7) in that the minimal  $N = \bar{N}$  iff (13) is an equality; this is so iff

$$(16) \quad \frac{1}{2}(\sigma_{(1)}^2) - \sigma_{(1)(k)} = \dots = \frac{1}{2}(\sigma_{(k-1)}^2) - \sigma_{(k-1)(k)},$$

which occurs if, e.g.,  $\sigma_1^2 = \dots = \sigma_k^2$  and  $\sigma_{(1)(k)} = \dots = \sigma_{(k-1)(k)}$ . For the latter to be satisfied it is sufficient (but not necessary) that *all* covariances be equal. (Note that the case when all covariances are equal and they and  $\sigma$  are known has been mentioned by Milton (1963), pp. 5-6, whose tables may then be used.)

$\bar{N}$  will be "exact" and approximation (15) identical to approximation (7) in the case  $\sigma_1^2 = \dots = \sigma_k^2 = \sigma^2$  and all covariances are equal to 0 (the case of independence). Now, if one knew that  $\sigma_1^2 = \dots = \sigma_k^2 = \sigma^2$  and that components of any  $k$ -vector tended to vary together if at all (i.e.  $\sigma_{ij} \geq 0$  ( $i, j = 1, \dots, k$ )), then one might extend the LFC and choose the "worst" such covariances (those which would maximize the approximation to the "exact"  $\bar{N}$ ). One sees (from criterion (16) and the role of  $\sigma_M$  in (13)) that all covariances equal to 0 accomplish this.

One may consider such extended LFC's in other cases specified by one's knowledge. If one knows  $\sigma_1^2 = \dots = \sigma_k^2 = \sigma^2$  and  $\sigma_{ij} = C$  ( $i \neq j; i, j = 1, \dots, k$ ) then  $\bar{N}$  will be "exact" and the worst value of  $C$  is the smallest possible (smallest such that the variance-covariance matrix is positive definite). Since for this case the matrix is positive definite iff  $\sigma^2 > C > -\sigma^2/(k - 1)$ , as  $C \downarrow -\sigma^2/(k - 1), \sigma_M \uparrow [k/(k - 1)]^{\frac{1}{2}}\sigma$ . Thus, in this case approximation (15) to  $\bar{N}$  will differ from approximation (7) by a factor  $k/(k - 1) = 1 + (k - 1)^{-1}$ .

Finally, since  $\sigma_M \leq 2^{\frac{1}{2}}\sigma$  (in the case  $\sigma_1^2 = \dots = \sigma_k^2 = \sigma^2$ ), the approximation to  $\bar{N}$  is at most twice that of (7). Thus, for large  $\alpha$  at *most* twice as many observations per population will be required as in the case of independence.

**5. Remarks.** We make the following remarks.

1. For the case  $k = 2$ , result (7) may be derived in a manner similar to that used by Dudewicz (1966), Section 2-B (i), in another context.
2. As the problem statement of Section 2 should intimate, result (7) is con-

jectured to hold for a general location parameter family whose observations obey the (classical) Central Limit Theorem. A similar result can be conjectured for Section 4, and is not intimidated by our problem statement there only so as to avoid a cumbersome notation.

3. In their monograph, Bechhofer, Kiefer, and Sobel (1968) have independently previously obtained (7) via a different proof. Although they obtain higher order corrections, their proof does not appear to generalize to multivariate problems. (See their Theorem 6.2.1, eq. 6.4.1, eq. 6.4.3, Lemma 6.5.1, and eq. 14.2.10.)

4. Thanks are especially due to Professor Robert E. Bechhofer for his suggestions for and guidance of Dudewicz (1966), which contains result (7), and to a referee for substantial simplification of the presentation and for the reference to Gupta (1963).

#### REFERENCES

- [1] BECHHOFER, R. E. (1954). A single-sample multiple decision procedure for ranking means of normal populations with known variances. *Ann. Math. Statist.* **25** 16-39.
- [2] BECHHOFER, R. E., KIEFER, J. and SOBEL, M. (1968). *Sequential Identification and Ranking Procedures (with special reference to Koopman-Darmonis populations)*. Univ. of Chicago Press.
- [3] DUDEWICZ, E. J. (1966). The efficiency of a nonparametric selection procedure: largest location parameter case. Technical Report No. 14, Department of Operations Research, Cornell Univ.
- [4] FELLER, W. (1957). *An Introduction to Probability Theory and Its Applications, Volume I* (Second Edition). Wiley, New York.
- [5] GUPTA, S. S. (1963). Probability integrals of multivariate normal and multivariate  $t$ . *Ann. Math. Statist.* **34** 792-828.
- [6] MILTON, R. C. (1963). Tables of the equally correlated multivariate normal probability integral. Technical Report No. 27, Department of Statistics, Univ. of Minnesota.
- [7] TEICHROEW, D. (1955). Probabilities Associated with Order Statistics in Samples from Two Normal Populations with Equal Variance. Chemical Corps Engineering Agency, Engineering Statistics Unit, Army Chemical Center, Maryland, December 7, ENASR No. ES-3.