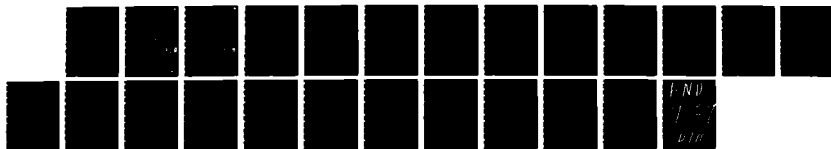
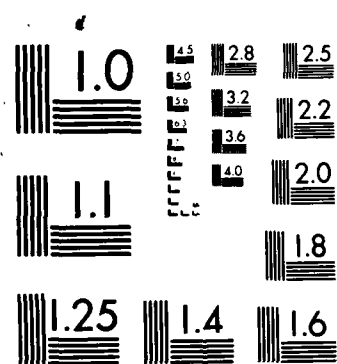


AD-A181 866 ON A LOWER CONFIDENCE BOUND FOR THE PROBABILITY OF A 171
CORRECT SELECTION: A (U) PURDUE UNIV LAFAYETTE IN DEPT
OF STATISTICS S S GUPTA ET AL MAR 87 TR-87-5
UNCLASSIFIED N88014-84-C-0167 F/G 12/3 NL





AD-A181 066

DTIC FILE COPY

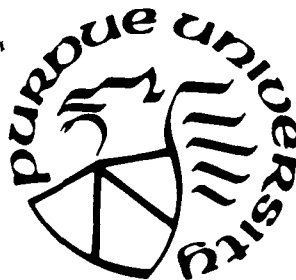
ON A LOWER CONFIDENCE BOUND FOR THE
PROBABILITY OF A CORRECT SELECTION: ANALYTICAL
AND SIMULATION STUDIES*

by

Shanti S. Gupta
Purdue University

TaChen Liang
Purdue University
Technical Report #87-5

PURDUE UNIVERSITY



DEPARTMENT OF STATISTICS

DISTRIBUTION STATEMENT A

Approved for public release
Distribution Unlimited

07 5 26 001
~~001~~

ON A LOWER CONFIDENCE BOUND FOR THE
PROBABILITY OF A CORRECT SELECTION: ANALYTICAL
AND SIMULATION STUDIES*

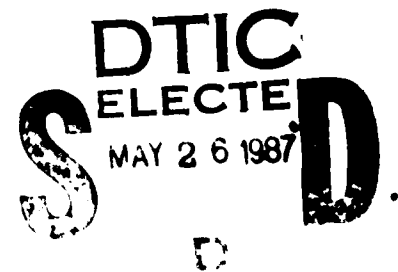
by

Shanti S. Gupta
Purdue University

TaChen Liang
Purdue University
Technical Report#87-5

Department of Statistics
Purdue University

March 1987



DISTRIBUTION STATEMENT A

Approved for public release
Distribution Unlimited

* The research of this author was supported in part by the Office of Naval Research Contract N00014-84-C-0167 at Purdue University and NSF Grant DMS-8606964. Reproduction in whole or in part is permitted for any purpose of the United States Government.

ON A LOWER CONFIDENCE BOUND FOR THE
PROBABILITY OF A CORRECT SELECTION: ANALYTICAL
AND SIMULATION STUDIES

by

Shanti S. Gupta* and TaChen Liang
Purdue University

Abstract

For the problem of selecting the best of several populations using the indifference (preference) zone formulation, a natural rule is to select the population yielding the largest sample value of an appropriate statistic. For this approach, it is required that the experimenter specify a number δ^* , say, which is a lower bound on the difference (separation) between the largest and the second largest parameter. However, in many real situations, it is hard to assign the value of δ^* and, therefore, in case that the assumption of indifference zone is violated, the probability of a correct selection cannot be guaranteed to be at least P^* , a prespecified value. In this paper, we are concerned with deriving a lower confidence bound for the probability of a correct selection for the general location model $F(x - \theta_i), i = 1, \dots, k$. First, we derive simultaneous lower confidence bounds on the differences between the largest (best) and each of the other non-best population parameters. Based on these, we obtain a lower confidence bound for the probability of a correct selection. The general result is then applied to the selection of the best mean of k normal populations with both the known and unknown common variances. In the first case one needs a single-stage procedure while in the second case a two-stage procedure is required. Some simulation investigations are described and their results are provided.

* Invited paper to be presented by this author at the First International Conference on Statistical Computing, 30 March - 2 April, 1987, to be held at Cesme, Izmir, Turkey.

This research was partially supported by the Office of Naval Research Contract N00014 84-0167 and NSF Grant DMS 8606964 at Purdue University



1
[]
[]

Copies
for
at

A-1

AMS 1980 Subject Classification: Primary 62F25, 62F07

KEY WORDS: Correct selection, Indifference zone, Lower confidence bound, Best population

1. Introduction

Let X_{ij} , $j = 1, \dots, n$, be n independent observations from a population π_i , where $\pi_1, \pi_2, \dots, \pi_k$ are independently distributed with continuous cumulative distribution function $G(x - \theta_i)$, $1 \leq i \leq k$, respectively. Let $\underline{\theta} = (\theta_1, \dots, \theta_k)$ and let $\theta_{(1)} \leq \dots \leq \theta_{(k)}$ denote the ordered values of $\theta_1, \dots, \theta_k$. It is assumed that the exact pairing between the ordered parameters and the unordered parameters is unknown. The population associated with the largest location parameter $\theta_{(k)}$ is called the best population. Assume that the experimenter is interested in the selection of the best population. For this purpose, we choose an appropriate statistic $Y_i = Y(X_{i1}, \dots, X_{in})$ with cumulative distribution function $F_n(y - \theta_i)$ and use the natural selection rule that selects the population yielding the largest Y_i as the best population. Let CS (correct selection) denote the event that the best population is selected. Then, the probability of a correct selection (PCS) applying the natural selection rule is

$$P_{\underline{\theta}}\{CS\} = \int_{-\infty}^{\infty} \prod_{i=1}^{k-1} F_n(y + \theta_{(k)} - \theta_{(i)}) dF_n(y). \quad (1.1)$$

To guarantee the probability of a correct selection, Bechhofer (1954) introduced the indifference zone approach in which the experimenter is asked to assign a positive value δ^* such that

$$\theta_{(k)} \geq \theta_{(k-1)} + \delta^*. \quad (1.2)$$

Thus, the subspace $\Omega(\delta^*) = \{\underline{\theta} | \theta_{(k)} \geq \theta_{(k-1)} + \delta^*\}$ is called the preference zone and its complement $\Omega^c(\delta^*) = \{\underline{\theta} | \theta_{(k)} < \theta_{(k-1)} + \delta^*\}$ is the indifference zone. We also let $\Omega = \Omega(\delta^*) \cup \Omega^c(\delta^*)$. On $\Omega(\delta^*)$, we have,

$$\inf_{\underline{\theta} \in \Omega(\delta^*)} P_{\underline{\theta}}\{CS\} = \int_{-\infty}^{\infty} [F_n(y + \delta^*)]^{k-1} dF_n(y) \quad (1.3)$$

Suppose that the function on the right-hand-side of (1.3) is an increasing function of the common sample size n and tends to one as n tends to infinity. Then, for a given probability $P^*(k-1) < P^* < 1$,

1), the minimum common sample size n_0 that is required to guarantee the probability of a correct selection to be at least P^* over the preference zone is determined by

$$n_0 \equiv n_0(\delta^*, P^*) = \min\{n \mid \int_{-\infty}^{\infty} [F_n(y + \delta^*)]^{k-1} dF_n(y) \geq P^*\}. \quad (1.4)$$

However, in a real situation, it may be hard to assign the value of δ^* such that $\theta_{(k)} \geq \theta_{(k-1)} + \delta^*$ since the parameter values $\theta_{(k)}, \theta_{(k-1)}$ are unknown. So that if the above assumption is not satisfied, then the probability of a correct selection cannot be guaranteed to be at least equal to P^* .

Recently, retrospective analyses regarding the PCS have been studied by some authors. Olkin, Sobel and Tong (1976, 1982) and Gibbons, Olkin and Sobel (1977) have presented estimators of the PCS. Faltin and McCulloch (1983) have studied the small-sample properties of the Olkin-Sobel-Tong's estimator of the PCS for the case when $k = 2$. Anderson, Bishop and Dudewicz (1977) gave a lower confidence bound on the PCS in the case of normal populations having a common variance which is either known or unknown. Kim (1986) presented a lower confidence bound on the PCS for the case where the underlying probability density function $f_n(y - \theta)$ of $F_n(y - \theta)$ has the monotone likelihood ratio property in y and θ and studied its application to the case of normal populations with common known or common unknown variances.

In this paper we are concerned with deriving a lower confidence bound for the probability of a correct selection for the general location model $G(x - \theta_i)$, $i = 1, \dots, k$. First, we derive simultaneous lower confidence bounds on the differences between the largest (best) and each of the other non-best population parameters. Based on these, we obtain a lower confidence bound for the probability of a correct selection. The general result is then applied to the selection of the best mean of k normal populations with both the known and unknown common variances. In the first case one needs a single-stage procedure while in the second case a two-stage procedure is required. Some simulation investigations are described and their results are provided.

2. A Lower Confidence Bound on PCS

For given δ^* and P^* , let n_0 be the minimum common sample size determined by (1.4). Let $Y_i = Y(X_{i1}, \dots, X_{in_0})$, be an appropriate statistic for inference regarding θ_i and let us assume that the distribution of $Y_i - \theta_i$ is independent of θ_i , $1 \leq i \leq k$. Let $Y_{[1]} \leq \dots \leq Y_{[k]}$ denote the order statistics of Y_i , $1 \leq i \leq k$. Also, let $\theta_{[i]}$ denote the (unknown) parameter associated with $Y_{[i]}$. For given α , $0 < \alpha < 1$, let $c(k, n_0, \alpha)$ be the value such that

$$P_{\theta} \left\{ \max_{1 \leq i \leq k} (Y_i - \theta_i) - \min_{1 \leq j \leq k} (Y_j - \theta_j) \leq c(k, n_0, \alpha) \right\} = 1 - \alpha. \quad (2.1)$$

Let $E = \left\{ \max_{1 \leq i \leq k} (Y_i - \theta_i) - \min_{1 \leq j \leq k} (Y_j - \theta_j) \leq c(k, n_0, \alpha) \right\}$. Then, we have the following lemma.

Lemma 2.1. $E \subset \{(Y_{[k]} - Y_{[i]} - c(k, n_0, \alpha))^+ \leq \theta_{(k)} - \theta_{(i)}, 1 \leq i \leq k-1\}$, where $(y)^+ = \max(0, y)$.

Proof: First note that for each $i = 1, \dots, k$,

$$\begin{aligned} \min_{j \leq i} (Y_{[j]} - \theta_{[j]}) &\leq \min_{j \leq i} (Y_{[i]} - \theta_{[j]}) \\ &= Y_{[i]} - \max_{j \leq i} \theta_{[j]} \\ &\leq Y_{[i]} - \theta_{(i)}. \end{aligned} \quad (2.2)$$

Thus,

$$\begin{aligned} E &\equiv \left\{ \max_{1 \leq i \leq k} (Y_i - \theta_i) - \min_{1 \leq j \leq k} (Y_j - \theta_j) \leq c(k, n_0, \alpha) \right\} \\ &\subset \left\{ \max_{1 \leq i \leq k} (Y_i - \theta_{(k)}) - \min_{1 \leq j \leq k} (Y_{[j]} - \theta_{[j]}) \leq c(k, n_0, \alpha) \right\} \\ &\subset \left\{ (Y_{[k]} - \theta_{(k)}) - \min_{1 \leq j \leq k-1} (Y_{[j]} - \theta_{[j]}) \leq c(k, n_0, \alpha) \right\} \\ &= \left\{ (Y_{[k]} - \theta_{(k)}) - \min_{j \leq i} (Y_{[i]} - \theta_{[j]}) \leq c(k, n_0, \alpha), 1 \leq i \leq k-1 \right\} \\ &\subset \left\{ (Y_{[k]} - \theta_{(k)}) - (Y_{[i]} - \theta_{(i)}) \leq c(k, n_0, \alpha), 1 \leq i \leq k-1 \right\} \text{ (by (2.2))} \\ &= \left\{ Y_{[k]} - Y_{[i]} - c(k, n_0, \alpha) \leq \theta_{(k)} - \theta_{(i)}, 1 \leq i \leq k-1 \right\} \\ &= \left\{ (Y_{[k]} - Y_{[i]} - c(k, n_0, \alpha))^+ \leq \theta_{(k)} - \theta_{(i)}, 1 \leq i \leq k-1 \right\}. \end{aligned}$$

Note that the last equality follows from the fact that $\theta_{(k)} - \theta_{(i)} \geq 0$ for all $i \leq k-1$. Hence, we complete the proof of this lemma.

Note that in (1.1), the probability of a correct selection $P_{\underline{\theta}}\{CS\}$ depends on the parameters $\underline{\theta} = (\theta_1, \dots, \theta_k)$ only via the differences $\theta_{(k)} - \theta_{(i)}$, $1 \leq i \leq k-1$. For convenience, we write $P_{\underline{\theta}}\{CS\} = P(\delta_1, \dots, \delta_{k-1})$ where $\delta_i = \theta_{(k)} - \theta_{(i)}$, $1 \leq i \leq k-1$. We see that $P(\delta_1, \dots, \delta_{k-1})$ is a nondecreasing function of δ_i for each $i = 1, 2, \dots, k-1$.

For each $i = 1, \dots, k-1$, let

$$\hat{\delta}_{L,i} = (Y_{[k]} - Y_{[i]} - c(k, n_0, \alpha))^+, \quad (2.3)$$

$$\hat{P}_L = P(\hat{\delta}_{L,1}, \dots, \hat{\delta}_{L,k-1}). \quad (2.4)$$

We propose \hat{P}_L as an estimator of a lower bound of the PCS. We have the following theorem.

Theorem 2.2. $P_{\underline{\theta}}\{P_{\underline{\theta}}\{CS\} \geq \hat{P}_L\} \geq 1 - \alpha$ for all $\underline{\theta} \in \Omega$.

Proof: By nondecreasing property of $P(\delta_1, \dots, \delta_{k-1})$ with respect to δ_i , $1 \leq i \leq k-1$, from (2.1) and Lemma 2.1, we have, for $\underline{\theta} \in \Omega$,

$$\begin{aligned} 1 - \alpha &= P_{\underline{\theta}}\{E\} \\ &\leq P_{\underline{\theta}}\{\hat{\delta}_{L,i} \leq \theta_{(k)} - \theta_{(i)}, 1 \leq i \leq k-1\} \\ &\leq P_{\underline{\theta}}\{P(\hat{\delta}_{L,1}, \dots, \hat{\delta}_{L,k-1}) \leq P(\delta_1, \dots, \delta_{k-1})\} \\ &= P_{\underline{\theta}}\{\hat{P}_L \leq P_{\underline{\theta}}\{CS\}\}. \end{aligned}$$

This completes the proof of this theorem.

3. Selection of the Best Normal Population in Terms of Means

Let X_{ij} , $1 \leq j \leq n$ be independent observations from $N(\theta_i, \sigma^2)$, $i = 1, \dots, k$ where the common variance σ^2 may be either known or unknown. The best population is the one associated

with the largest mean $\theta_{(k)}$. We consider two situations according to whether the common variance σ^2 is known or unknown.

3.1. Lower Confidence Bound for PCS : σ^2 Known Case.

When the value of the common variance σ^2 is known, for $\theta \in \Omega$, the probability of a correct selection applying the natural selection rule is:

$$P_{\theta}\{CS\} = \int_{-\infty}^{\infty} \prod_{i=1}^{k-1} \Phi\left(x + \frac{\sqrt{n_0}(\theta_{(k)} - \theta_{(i)})}{\sigma}\right) d\Phi(x), \quad (3.1)$$

where $\Phi(\cdot)$ is the standard normal distribution function, and the value of the sample size n_0 , for the indifference zone formulation, is determined by

$$n_0 = \min\{n \mid \int_{-\infty}^{\infty} [\Phi(x + \frac{\sqrt{n}\delta^*}{\sigma})]^{k-1} d\Phi(x) \geq P^*\}. \quad (3.2)$$

Let $\bar{X}_i = \frac{1}{n_0} \sum_{j=1}^{n_0} X_{ij}$. For given $0 < \alpha < 1$, choose the value $c(k, n_0, \alpha)$ such that

$$P_{\theta}\left\{\max_{1 \leq i \leq k} (\bar{X}_i - \theta_i) - \min_{1 \leq j \leq k} (X_j - \theta_j) \leq c(k, n_0, \alpha)\right\} = 1 - \alpha. \quad (3.3)$$

Note that here, $c(k, n_0, \alpha) = \frac{\sigma}{\sqrt{n_0}} q_{k, \infty}^{\alpha}$, where $q_{k, \infty}^{\alpha}$ is the $100(1 - \alpha)\%$ th percentile of Tukey's studentized range statistic with parameters (k, ∞) . The value of $q_{k, \infty}^{\alpha}$ is available from Harter (1969). Then, we define

$$\hat{\delta}_{L,i} = (\bar{X}_{[k]} - \bar{X}_{[i]} - c(k, n_0, \alpha))^+ \quad (3.4)$$

and

$$\hat{P}_L = P(\hat{\delta}_{L,1}, \dots, \hat{\delta}_{L,k-1}) = \int_{-\infty}^{\infty} \prod_{i=1}^{k-1} \Phi\left(x + \frac{\sqrt{n_0}\hat{\delta}_{L,i}}{\sigma}\right) d\Phi(x). \quad (3.5)$$

Then, by Theorem 2.2, $P_{\theta}\{P_{\theta}\{CS\} \geq \hat{P}_L\} \geq 1 - \alpha$ for all $\theta \in \Omega$.

3.2. Lower Confidence Bound for PCS : σ^2 Unknown Case.

When the common variance σ^2 is unknown, Bechhofer, Dunnett and Sobel (1954) presented a two-stage selection rule, which is briefly described as follows.

Take a first sample of n_0 ($n_0 \geq 2$) observations from each of the k populations. Compute $\bar{X}_i = \frac{1}{n_0} \sum_{j=1}^{n_0} X_{ij}$, ($1 \leq i \leq k$), and $S^2 = \frac{1}{k(n_0-1)} \sum_{i=1}^k \sum_{j=1}^{n_0} (X_{ij} - \bar{X}_i)^2$. Define $N = \max\{n_0, \lceil \frac{S^2 h^2}{\delta^*} \rceil\}$ where the symbol $\lceil y \rceil$ denotes the smallest integer not less than y , and h is a positive value such that

$$\int_0^\infty \int_{-\infty}^\infty [\Phi(x + wh)]^{k-1} d\Phi(x) dF_W(w) = P^*, \quad (3.6)$$

$k^{-1} < P^* < 1$, and $F_W(\cdot)$ is the distribution function of the nonnegative random variable W with $k(n_0 - 1)W^2$ following $\chi^2(k(n_0 - 1))$ distribution.

Then, take additional $N - n_0$ observations from each population. Compute the overall mean $\bar{X}_i(N) = \frac{1}{N} \sum_{j=1}^N X_{ij}$, $1 \leq i \leq k$. We then select the population yielding the largest observation $\bar{X}_{[k]}(N)$ as the best population.

For this two-stage selection rule, the probability of a correct selection is:

$$\begin{aligned} P_{\theta}\{CS\} &= P_{\theta}\{\bar{X}_{(k)}(N) > \bar{X}_{(i)}(N), i \neq k\} \\ &= P_{\theta}\left\{\frac{\sqrt{N}(\bar{X}_{(k)}(N) - \theta_{(k)})}{\sigma} + \frac{\sqrt{N}(\theta_{(k)} - \theta_{(i)})}{\sigma} > \frac{\sqrt{N}(\bar{X}_{(i)} - \theta_{(i)})}{\sigma}, i \neq k\right\} \\ &\geq P_{\theta}\left\{\frac{\sqrt{N}(\bar{X}_{(k)}(N) - \theta_{(k)})}{\sigma} + \frac{h(\theta_{(k)} - \theta_{(i)})}{\delta^*} \frac{S}{\sigma} > \frac{\sqrt{N}(\bar{X}_{(i)} - \theta_{(i)})}{\sigma}, i \neq k\right\} \\ &\quad \left(\text{since } N \geq \left\lceil \frac{S^2 h^2}{\delta^*} \right\rceil\right) \\ &= P\left\{Z_k + \frac{h(\theta_{(k)} - \theta_{(i)})}{\delta^*} W > Z_i, i \neq k\right\} \\ &= \int_0^\infty \int_{-\infty}^\infty \prod_{i=1}^{k-1} \Phi\left(z + \frac{h(\theta_{(k)} - \theta_{(i)})}{\delta^*} w\right) d\Phi(z) dF_W(w), \end{aligned} \quad (3.7)$$

where Z_1, \dots, Z_k are iid random variables having standard normal distribution, and $W = S/\sigma$ with $k(n_0 - 1)W^2 \sim \chi^2(k(n_0 - 1))$ and (Z_1, \dots, Z_k) and W are independent.

Thus, to obtain a lower confidence bound for $P_{\underline{\theta}}\{CS\}$, it suffices to find simultaneous lower confidence bounds for $\theta_{(k)} - \theta_{(i)}$, $1 \leq i \leq k-1$. Then, replacing the $\theta_{(k)} - \theta_{(i)}$, $1 \leq i \leq k-1$, by the corresponding lower confidence bounds into the function on the right-hand-side of (3.7), we obtain a lower confidence bound for $P_{\underline{\theta}}\{CS\}$. For convenience, we let

$$Q(\delta_1, \dots, \delta_{k-1}) = \int_0^\infty \int_{-\infty}^\infty \prod_{i=1}^{k-1} \Phi\left(z + \frac{h(\theta_{(k)} - \theta_{(i)})w}{\delta^*}\right) d\Phi(z) dF_W(w). \quad (3.8)$$

Let $c = Sq_{k,k(n_0-1)}^\alpha / \sqrt{N}$, where $q_{k,k(n_0-1)}^\alpha$ is the $100(1-\alpha)\%$ th percentile of Tukey's studentized range statistic with parameters $(k, k(n_0-1))$. Define

$$\hat{\delta}_{L,i} = (\bar{X}_{[k]}(N) - \bar{X}_{[i]}(N) - c)^+, \quad (3.9)$$

and

$$\hat{Q}_L = Q(\hat{\delta}_{L,1}, \dots, \hat{\delta}_{L,k-1}). \quad (3.10)$$

We propose \hat{Q}_L as an estimator of a lower bound of $P_{\underline{\theta}}\{CS\}$.

Lemma 3.1. Let $E = \{\max_{1 \leq i \leq k} (\bar{X}_i(N) - \theta_i) - \min_{1 \leq j \leq k} (\bar{X}_j(N) - \theta_j) \leq c\}$. Then, $P_{\underline{\theta}}\{E\} = 1 - \alpha$ for all $\underline{\theta} \in \Omega$.

Proof:

$$\begin{aligned} P_{\underline{\theta}}(E) &= P_{\underline{\theta}}\{\max_{1 \leq i \leq k} (\bar{X}_i(N) - \theta_i) - \min_{1 \leq j \leq k} (\bar{X}_j(N) - \theta_j) \leq c\} \\ &= P_{\underline{\theta}}\{\max_{1 \leq i \leq k} \sqrt{N}(\bar{X}_i(N) - \theta_i) - \min_{1 \leq j \leq k} \sqrt{N}(\bar{X}_j(N) - \theta_j) \leq Sq_{k,k(n_0-1)}^\alpha\} \\ &= 1 - \alpha, \end{aligned}$$

where the last equality follows from the definition of $q_{k,k(n_0-1)}^\alpha$.

Lemma 3.2. $P_{\underline{\theta}}\{\hat{\delta}_{L,i} \leq \theta_{(k)} - \theta_{(i)}, 1 \leq i \leq k-1\} \geq 1 - \alpha$ for all $\underline{\theta} \in \Omega$.

Proof: Following the same argument as in Lemma 2.1, we have $E \subset \{\hat{\delta}_{L,i} \leq \theta_{(k)} - \theta_{(i)}, 1 \leq i \leq k-1\}$. Then using Lemma 3.1 leads to the conclusion of Lemma 3.2.

Lemma 3.2 and the increasing property of the function $Q(\delta_1, \dots, \delta_{k-1})$ with respect to $\delta_i, 1 \leq i \leq k-1$, lead to the following main result.

Theorem 3.3. $P_{\theta}\{P_{\theta}\{CS\} \geq \hat{Q}_L\} \geq 1 - \alpha$ for all $\theta \in \Omega$.

Proof: Note that $P_{\theta}\{CS\} \geq Q(\delta_1, \dots, \delta_{k-1})$ for all $\theta \in \Omega$. Therefore, $P_{\theta}\{P_{\theta}\{CS\} \geq \hat{Q}_L\} \geq P_{\theta}\{Q(\delta_1, \dots, \delta_{k-1}) \geq \hat{Q}_L\} \geq 1 - \alpha$ for all $\theta \in \Omega$.

4. Remark and Example

Anderson, Bishop and Dudewicz (1977) and Kim (1986) have also studied the problem of finding a lower confidence bound on PCS. They considered the retrospective analysis to approach a lower confidence bound for PCS no matter what the sampling rule is. However, our approach is different from theirs. We use the following example to illustrate our procedure and describe the difference between ours and Kim's approach.

Example (The data is taken from Problem 3.1, page 97, of Gibbons, Olkin and Sobel (1977)).

The experimenter wants to compare dry shear strength of $k = 6$ different resin glues for bonding yellow birch plywood. Assume that the distributions of the strength for each glue are normal with common unknown variance σ^2 . Based on some past information, the experimenter assigns $\delta^* = 20$. Then, using indifference zone formulation, a two-stage natural selection rule is applied here. Let $P^* = 0.90$ and let the initial sample size of n_0 be 6. The observations (readings) are taken to measure the strength of the glue. Thus, large values are more desirable in this application. The data are given in the upper part of Table 1.

Now, $N = \max\{n_0, \lceil \frac{S^2 h^2}{\delta^*} \rceil\}$. For $k = 6, n_0 = 6, P^* = 0.90$, from Gupta, Panchapakesan and Sohn (1985), $h = 1.97982\sqrt{2}$. Therefore, $N = 10$ and hence $N - n_0 = 4$

Table 1. Shear Strength of Six Types of Glue

	Glue					
	1	2	3	4	5	6
observations taken at the first-stage	102	70	100	120	151	220
	58	83	102	110	156	243
	45	78	80	182	192	189
	79	93	119	130	162	176
	68	98	59	95	166	176
	63	66	99	143	158	181
n_0	6	6	6	6	6	6
\bar{X}_i	69.17	81.33	93.17	130.00	164.17	197.50
$S^2 = \frac{1}{k(n_0-1)} \sum_{i=1}^k \sum_{j=1}^{n_0} (X_{ij} - \bar{X}_i)^2 = 479.51$						
observations taken at the second-stage	117	92	100	113	173	206
	94	79	109	140	157	233
	99	134	128	123	233	162
	63	131	138	132	238	179
N	10	10	10	10	10	10
$\bar{X}_i(N)$	78.8	92.4	103.4	128.8	178.6	196.5
$S^2(N) = \frac{1}{k(N-1)} \sum_{i=1}^k \sum_{j=1}^N (X_{ij} - \bar{X}_i(N))^2 = 656.98$						

additional observations should be taken from each population. The observations taken at the second-stage are given in the lower part of Table 1.

We then have the overall sample means: $\bar{X}_1(N) = 78.8$, $\bar{X}_2(N) = 92.4$, $\bar{X}_3(N) = 103.4$, $\bar{X}_4(N) = 128.8$, $\bar{X}_5(N) = 178.6$, $\bar{X}_6(N) = 196.5$. According to the two-stage natural selection rule, Glue 6 which yields the largest sample mean is selected as the best.

However, we do not know whether the largest and the second largest unknown means differ at least by $\delta^* = 20$ or not. A reasonable question is: What kind of confidence statement can be made regarding the PCS? By the method described in Section 3.2, for $\alpha = 0.10$, from Harter (1969), $q_{k, k(n_0-1)}^\alpha = 3.851$. Thus, $c = Sq_{k, k(n_0-1)}^\alpha / \sqrt{N} = 26.667$. Therefore, $\hat{\delta}_{L,1} = 91.033$, $\hat{\delta}_{L,2} =$

77.433, $\hat{\delta}_{L,3} = 66.433$, $\hat{\delta}_{L,4} = 41.033$, $\hat{\delta}_{L,5} = 0$. After some computation, we have $\hat{Q}_L = 0.5000$. Therefore, we can state that with at least 90% confidence that $PCS \geq \hat{Q}_L = 0.5000$ for all values of true unknown means.

For the problem of selecting the largest normal mean from among k ($k \geq 2$) normal populations having a common unknown variance σ^2 , Kim (1986) proposed a lower confidence bound using a retrospective analysis based on the observations X_{ij} , $1 \leq i \leq k$, $1 \leq j \leq n$, where the common sample size n is arbitrary. His proposed lower confidence bound is given as follows: With at least $100(1 - \alpha)\%$ confidence,

$$PCS \geq \int_{-\infty}^{\infty} \Phi^{k-1} \left[x + \sqrt{2} h_{\nu} \left(\frac{\sqrt{n}(\bar{X}_{[k]} - \bar{X}_{[k-1]})}{\sqrt{2}S} \right) \right] d\Phi(x) \quad (4.1)$$

where $\nu = k(n - 1)$ and the function $h_{\nu}(\cdot)$ is implicitly defined by

$$\int_0^{\infty} [\Phi(h_{\nu}(t) - tw) + \Phi(-h_{\nu}(t) - tw)] dF_W(w) = \alpha \quad (4.2)$$

for $t \geq t_{\frac{\alpha}{2}}(\nu)$ and $h_{\nu}(t) = 0$ for $0 \leq t \leq t_{\frac{\alpha}{2}}(\nu)$. Here, $t_{\frac{\alpha}{2}}(\nu)$ is the upper $\frac{\alpha}{2}$ quantile of the t distribution with ν degrees of freedom and $F_W(\cdot)$ is the distribution of a non-negative random variable W with $\nu W^2 \sim \chi^2(\nu)$. Kim (1986) also provided some tables for the $h_{\nu}(t)$ values to implement his procedure.

For the data set given above, following Kim's procedure, we have $n = 10$, $\nu = k(n - 1) = 54$, and the pooled sample standard deviation based on the total 60 observations is $S = 25.63$. Thus, $t = \frac{\sqrt{n}(\bar{X}_{[k]} - \bar{X}_{[k-1]})}{\sqrt{2}S} = 1.562 < t_{\frac{\alpha}{2}}(\nu) \approx 1.671$ where $\alpha = 0.10$. Therefore, by the definition of $h_{\nu}(\cdot)$, $h_{\nu}(\frac{\sqrt{n}(\bar{X}_{[k]} - \bar{X}_{[k-1]})}{\sqrt{2}S}) = 0$. Then, by (4.1), one can only claim: With at least 90% confidence, $PCS \geq \frac{1}{k} = \frac{1}{6}$. Clearly, our lower confidence bound, in this example, is better than that of Kim.

5. Simulation Studies

For the normal means selection problem, for various parameter configurations, the behaviors of \hat{P}_L and \hat{Q}_L were simulated. Two types of parameter configurations were simulated: a slippage

configuration $\theta_{(1)} = \dots = \theta_{(k-1)} = \theta_{(k)} - \Delta$ and an equally spaced configuration $\theta_{(i)} - \theta_{(i-1)} = \Delta$, $i = 2, \dots, k$. For simulation, we suppose that the assigned value of δ^* is 1 and also the assigned probability levels are $P^* = 0.90$ and $P^* = 0.95$. When the common variance σ^2 is known, the common sample size n_0 is determined by (3.2). When σ^2 is unknown, the initial common sample size is set equal to ten. The simulation process was repeated $M = 1000$ times for the case where σ^2 is known and $M = 400$ times for the σ^2 unknown case. For each simulation, the random observation X_{ij} is generated from $N(\theta_i, \sigma^2)$ with $\sigma^2 = 1$. The values of \hat{P}_L and \hat{Q}_L were computed. The averages of the 1000 \hat{P}_L and 400 \hat{Q}_L are reported in Table 2 and Table 3, respectively. In each table, the numbers in the parentheses are the standard errors of the corresponding estimators.

For convenience, we let $\hat{P}_L(P^*, \Delta, \alpha, T)$ and $\hat{Q}_L(P^*, \Delta, \alpha, T)$ denote the corresponding \hat{P}_L and \hat{Q}_L for given values of P^*, Δ, α and T , where T denotes the type of parameter configuration. The slippage configuration is denoted by S and the equally spaced configuration is denoted by ES .

The simulation results indicate the following:

1. Note that for fixed P^*, α and T , the PCS is a nondecreasing function of Δ . Therefore, it is reasonable to expect that both $\hat{P}_L(P^*, \Delta, \alpha, T)$ and $\hat{Q}_L(P^*, \Delta, \alpha, T)$ be nondecreasing in Δ . The simulation results indicate that this is so.
2. For fixed P^*, α and Δ , the PCS under equally spaced parameter configuration is larger than the PCS under the slippage configuration. The simulation results also indicate that this behavior holds. That is, from the simulation results, we find: $\hat{P}_L(P^*, \Delta, \alpha, ES) > \hat{P}_L(P^*, \Delta, \alpha, S)$ and $\hat{Q}_L(P^*, \Delta, \alpha, ES) > \hat{Q}_L(P^*, \Delta, \alpha, S)$.
3. For fixed values P^*, Δ and T , the simulation results indicate that $\hat{P}_L(P^*, \Delta, 0.2, T) > \hat{P}_L(P^*, \Delta, 0.1, T)$ and $\hat{Q}_L(P^*, \Delta, 0.2, T) > \hat{Q}_L(P^*, \Delta, 0.1, T)$. These results are as expected since $q_{k,\nu}^\alpha$ is nondecreasing in α for fixed k and ν .

4. For fixed Δ, α and T , $\hat{P}_L(P^*, \Delta, \alpha, T)$ is nondecreasing in P^* . Note that according to the sampling rule used in this paper, assigning large P^* -value implicitly implies taking more observations. Thus the simulation results seem to indicate that $\hat{P}_L(P^*, \Delta, \alpha, T)$ is nondecreasing in the sample size. For the σ^2 unknown case, for both $k = 3$ and 5, the corresponding values of $\hat{Q}_L(P^*, \Delta, \alpha, T)$ are also nondecreasing in P^* .

Table 2 Simulated Values of P_L for $k = 3$, σ^2 Known

Δ	90% lower confidence bound				80% lower confidence bound			
	Slippage		Equally Spaced		Slippage		Equally Spaced	
	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$
0.5	0.3497 (0.0016)	0.3573 (0.0020)	0.3708 (0.0023)	0.3993 (0.0032)	0.3648 (0.0023)	0.3779 (0.0029)	0.3975 (0.0031)	0.4340 (0.0040)
0.8	0.3763 (0.0026)	0.3989 (0.0036)	0.4389 (0.0037)	0.5037 (0.0047)	0.3995 (0.0035)	0.4400 (0.0047)	0.4846 (0.0045)	0.5553 (0.0053)
1.0	0.3960 (0.0035)	0.4485 (0.0049)	0.5007 (0.0045)	0.5823 (0.0052)	0.4376 (0.0045)	0.5057 (0.0058)	0.5555 (0.0052)	0.6395 (0.0057)
1.5	0.5140 (0.0057)	0.6463 (0.0067)	0.6610 (0.0057)	0.7679 (0.0057)	0.5553 (0.0064)	0.7229 (0.0066)	0.7231 (0.0058)	0.8254 (0.0053)
2.0	0.6876 (0.0064)	0.8539 (0.0052)	0.8660 (0.0054)	0.9151 (0.0039)	0.7635 (0.0062)	0.9036 (0.0042)	0.8597 (0.0049)	0.9458 (0.0040)

Table 2 (continued) Simulated Values of P_L for $k = 5$, σ^2 Known

Δ	90% lower confidence bound				80% lower confidence bound			
	Slippage		Equally Spaced		Slippage		Equally Spaced	
	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$
0.5	0.2575 (0.0007)	0.2163 (0.0011)	0.3054 (0.0025)	0.3459 (0.0020)	0.2158 (0.0011)	0.2204 (0.0016)	0.2398 (0.0032)	0.2814 (0.0030)
0.8	0.2237 (0.0016)	0.2382 (0.0025)	0.4371 (0.0038)	0.4874 (0.0036)	0.2419 (0.0023)	0.2651 (0.0023)	0.4786 (0.0045)	0.5310 (0.0046)
1.0	0.2460 (0.0023)	0.2798 (0.0037)	0.5133 (0.0049)	0.5690 (0.0047)	0.2774 (0.0025)	0.3243 (0.0048)	0.5693 (0.0051)	0.6190 (0.0053)
1.5	0.3813 (0.0050)	0.5093 (0.0069)	0.6820 (0.0057)	0.7685 (0.0057)	0.4536 (0.0065)	0.5839 (0.0072)	0.7381 (0.0058)	0.8191 (0.0055)
2.0	0.6234 (0.0072)	0.7907 (0.0061)	0.8430 (0.0057)	0.9228 (0.0038)	0.7076 (0.0060)	0.8541 (0.0051)	0.8855 (0.0046)	0.9489 (0.0030)

Table 2 (continued) Simulated Values of P_L for $k = 10$, σ^2 Known

Δ	90% lower confidence bound				80% lower confidence bound			
	Slippage		Equally Spaced		Slippage		Equally Spaced	
	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$
0.5	0.1021 (0.0003)	0.1029 (0.0003)	0.2906 (0.0023)	0.3260 (0.0026)	0.1045 (0.0004)	0.1059 (0.0005)	0.3146 (0.0026)	0.3534 (0.0030)
0.8	0.1090 (0.0007)	0.1144 (0.0016)	0.4185 (0.0032)	0.4668 (0.0036)	0.1164 (0.1111)	0.1251 (0.0016)	0.4526 (0.0037)	0.4981 (0.0041)
1.0	0.1213 (0.0014)	0.1365 (0.0021)	0.4984 (0.0038)	0.5434 (0.0042)	0.1355 (0.0020)	0.1588 (0.0029)	0.5345 (0.0044)	0.5798 (0.0048)
1.5	0.2187 (0.0045)	0.3118 (0.0061)	0.6643 (0.0056)	0.7403 (0.0058)	0.2709 (0.0055)	0.3843 (0.0069)	0.7096 (0.0059)	0.7879 (0.0057)
2.0	0.4773 (0.0076)	0.6713 (0.0073)	0.8397 (0.0052)	0.9167 (0.0039)	0.5647 (0.0078)	0.7508 (0.0067)	0.8797 (0.0046)	0.9423 (0.0032)

Table 3 Simulated Values of \hat{Q}_L for $k = 3$, σ^2 Unknown, $n_0 = 10$

Δ	90% lower confidence bound				80% lower confidence bound			
	Slippage		Equally Spaced		Slippage		Equally Spaced	
	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$
0.5	0.3467 (0.0021)	0.4307 (0.0037)	0.4141 (0.0043)	0.5158 (0.0037)	0.3689 (0.0037)	0.4757 (0.0082)	0.4633 (0.0005)	0.5791 (0.0058)
1.0	0.6047 (0.0034)	0.4831 (0.0063)	0.4813 (0.0056)	0.5933 (0.0040)	0.4931 (0.0055)	0.5759 (0.0081)	0.5352 (0.0064)	0.6641 (0.0076)
2.0	0.7068 (0.0082)	0.9466 (0.0024)	0.8203 (0.0006)	0.9664 (0.0014)	0.8054 (0.0070)	0.9737 (0.0014)	0.8869 (0.0054)	0.9839 (0.0007)

Table 3 (continued) Simulated Values of \hat{Q}_L for $k = 5$, σ^2 Unknown, $n_0 = 10$

Δ	90% lower confidence bound				80% lower confidence bound			
	Suprage		Equally Spaced		Suprage		Equally Spaced	
	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$	$P^* = 0.90$	$P^* = 0.95$
0.8	0.2066 (0.0006)	0.2887 (0.0014)	0.4391 (0.0024)	0.6999 (0.0024)	0.2306 (0.0014)	0.3347 (0.0020)	0.5251 (0.0036)	0.8113 (0.0024)
1.0	0.2331 (0.0015)	0.3492 (0.0022)	0.5564 (0.0034)	0.8757 (0.0013)	0.2790 (0.0024)	0.3963 (0.0029)	0.6732 (0.0042)	0.9328 (0.0008)
2.0	0.7844 (0.0050)	0.9430 (0.0017)	0.9717 (0.0015)	0.9994 (0.0003)	0.8701 (0.0040)	0.9717 (0.0011)	0.9860 (0.0013)	0.9998 (0.0000)

6. A Lower Confidence Bound on PCS for Scale Parameter Model

The results in Section 2 are derived for a location parameter model. Similar results for a scale parameter model can also be obtained. For the problem of selecting the population with the largest scale parameter $\theta_{(k)}$, the P^* in (1) is replaced by

$$P_{\theta} \{ \hat{S}_L \} = \int_0^{\infty} \prod_{i=1}^{k-1} F_{\theta_i} \left(\frac{y_i}{\theta_i} \right) dF_{\theta_k} \left(\frac{y}{\theta_k} \right), \quad \theta_k > 0, \quad (6.1)$$

where $F_{\theta_i}(y/\theta_i)$ is the cumulative distribution function of an appropriate nonnegative statistic Y_i , $Y_i(X_{i1}, \dots, X_{in_i})$, and the common sample size n is determined according to some sampling rule. Suppose that the distribution of Y_i/θ_i is independent of θ_i , $i = 1, \dots, k$. For given $\alpha \in (0, 1)$, let T be the smallest value such that

$$P_{\theta_k} \left(\max_{1 \leq i \leq k-1} Y_i/\theta_i \leq Y_k/\theta_k \mid T \right) = 1 - \alpha. \quad (6.2)$$

Note that $T \geq 1$ since $\max_{1 \leq i \leq k-1} Y_i/\theta_i \leq \max_{1 \leq i \leq k} Y_i/\theta_i \leq 1$.

Analogous to the result obtained in Section 2, we have

$$P_{\theta_k} \{ \theta_{(k)} \geq c_{\alpha} = Y_{(k)}/Y_{(k)} T \} = 1 - \alpha, \quad \theta_k > 0, \quad (6.3)$$

for all θ_k , where now the parameter θ_k is $\theta_{(k)}$, $\theta_1, \dots, \theta_{k-1} = \theta_{(k)}/T$, $i = 1, \dots, k-1$.

Let

$$T_{\alpha} = \inf \{ T \mid P_{\theta_k} \{ \max_{1 \leq i \leq k-1} Y_i/\theta_i \leq Y_k/\theta_k \mid T \} = 1 - \alpha \}. \quad (6.4)$$

where $(y)^0 = \max(y, 1)$. Replacing $\theta_{(k)}/\theta_{(1)}$ in (6.1) by $\delta_{L,1}$, we obtain

$$\hat{P}_L = \int_0^\infty \prod_{i=1}^{k-1} F_n(\delta_{L,i}y) dF_n(y). \quad (6.5)$$

We propose \hat{P}_L as an estimator of a lower bound for the PCS. We have

$$P_\theta\{P_\theta\{CS\} \geq \hat{P}_L\} \geq 1 - \alpha \text{ for all } \theta \in \Omega. \quad (6.6)$$

References

- Anderson, P. O., Bishop, T. A. and Dudewicz, E. J. (1977). Indifference zone ranking and selection confidence intervals for true achieved P(CD). *Commun. Statist.* **A(6)**, 1121-1132.
- 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.
- Bechhofer, R. E. and Dunnett, C. W. and Sobel, M. (1954). A two-sample multiple decision procedure for ranking means of normal populations with a common unknown variance. *Biometrika* **41**, 170-176.
- Edging, E. and McCulloch, C. E. (1983). On the small-sample properties of the Olkin-Sobel-Tong estimator of the probability of correct selection. *J. Amer. Statist. Assoc.* **78**, 464-467.
- Gilbons, J. D., Olkin, I. and Sobel, M. (1977). *Selecting and Ordering Populations: A New Statistical Methodology*. New York: John Wiley.
- Gupta, S. S., Panchapakesan, S. and Sobel, J. B. (1985). On the distribution of the studentized maximum of equally correlated normal random variables. *Commun. Statist. Simula. Computa.* **14(1)**, 1-3, 135.
- Harter, H. T. (1966). *Order Statistics: Their Theory and Application*. Vol. 1, *Tests Based on Range and Student's t-Range of Samples from a Normal Population*. Aerospace Research Laboratories.

- Kim, W.-C. (1986). A lower confidence bound on the probability of a correct selection. *J. Amer. Statist. Assoc.* **81**, 1012-1017.
- Olkin, I., Sobel, M. and Tong, Y. L. (1976). Estimating the true probability of a correct selection for location and scale parameter families. Technical Report 110, Stanford University, Department of Statistics.
- Olkin, I., Sobel, M. and Tong, Y. L. (1982). Bounds for a k -fold integral for location and scale parameter models with applications to statistical ranking and selection problem. *Statistical Decision Theory and Related Topics III* Vol. 2 (Eds. S. S. Gupta and J. O. Berger), New York, Academic Press, pp. 193-212.

BEFORE COMPLETING FORM		
1. REPORT NUMBER Technical Report #87-5	2. GOVT ACCESSION NO ADA181066	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) ON A LOWER CONFIDENCE BOUND FOR THE PROBABILITY OF A CORRECT SELECTION: ANALYTICAL AND SIMULATION STUDIES		5. TYPE OF REPORT & PERIOD COVERED Technical
7. AUTHOR(s) Shanti S. Gupta and TaChen Liang		6. PERFORMING ORG. REPORT NUMBER Technical Report #87-5
9. PERFORMING ORGANIZATION NAME AND ADDRESS Purdue University Department of Statistics West Lafayette, IN 47907		8. CONTRACT OR GRANT NUMBER(s) N00014-84-C-0167 and NSF Grant DMS-8606964
11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research Washington, DC		10. PROGRAM ELEMENT, PROJECT, TASK, AREA & WORK UNIT NUMBERS
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		12. REPORT DATE March 1987
		13. NUMBER OF PAGES 17
		15. SECURITY CLASS. (of this report) Unclassified
		16. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release, distribution unlimited.		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Correct selection, Indifference zone, Lower confidence bound, Best population		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) For the problem of selecting the best of several populations using the indifference (preference) zone formulation, a natural rule is to select the population yielding the largest sample value of an appropriate statistic. For this approach, it is required that the experimenter specify a number ϵ^* , say, which is a lower bound on the difference (separation) between the largest and the second largest parameter. However, in many real situations, it is hard to assign the value of ϵ^* and, therefore, in case that the assumption of indifference zone is violated, the probability of a correct selection cannot be guaranteed to be at least ϵ^* , a prestpecified value.		

In this paper, we are concerned with deriving a lower confidence bound for the probability of a correct selection for the general location model $F(x-\theta_i)$, $i = 1, \dots, k$. First, we derive simultaneous lower confidence bounds on the differences between the largest (best) and each of the other non-best population parameters. Based on these, we obtain a lower confidence bound for the probability of a correct selection. The general result is then applied to the selection of the best mean of k normal populations with both the known and unknown common variances. In the first case one needs a single-stage procedure while in the second case a two-stage procedure is required. Some simulation investigations are described and their results are provided.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

END

7-87

DTIC