BAYESIAN SUBSET SELECTION  FOR ADDITIVE

AND LINEAR LOSS FUNCTIONS*

by

Klaus-J. Miescke

University of Mainz and Purdue University

Department of Statistics
Division of Mathematical Sciences
Mimeograph Series #78-23

October 1978

# BAYESIAN SUBSET SELECTION FOR ADDITIVE
## AND LINEAR LOSS FUNCTIONS*

by

## Klaus-J. Miescke

### University of Mainz and Purdue University

1. <u>Introduction and Summary</u>.

Let $\{P_\theta\}_{\theta \in \Omega}$, $\Omega \subseteq \mathbb{R}$ be a given family of probability distributions over

$(\mathbb{R}, \mathcal{B})$, where $\mathcal{B}$ denotes the corresponding Borel-sets, and let $\underline{X}_i = (X_{i1}, \ldots, X_{in})$,

$i = 1, \ldots, k$, be $k$ independent samples of common size $n$ from $k$ populations $\pi_i$ with

distributions $P_{\theta_i}$, $i = 1, \ldots, k$, where $\theta_1, \ldots, \theta_k \in \Omega$ are unknown. If our goal is

to select a random non-empty subset $S(\underline{X}_1, \ldots, \underline{X}_k)$ of populations which is "good"

(i.e. is associated with large $\theta$-values in a certain sense), then there are several

possible choices for a suitable criterion to make the term "good" more precise.

One of these - a Bayes criterion with additive and especially linear loss functions -

will be the topic of this paper.

For simplicity we assume that the rule $S$ take values in $G = \{s \mid s \subseteq \{1, \ldots, k\},$

$s \neq \phi\}$, where $S = s$ now means that all $\pi_i$ are selected for which $i \in s$, $i = 1, \ldots, k$.

Moreover for every $s \in G$ let $|s|$ denote the size of $s$.

For a given loss function $L: \Omega^k \times G \to \mathbb{R}$ and a given prior distribution $\tau$ over

$(\Omega^k, \mathcal{B}_k)$ (for the now random parameter vector $\underline{\theta}$ with values $\underline{\theta} \in \Omega^k$) every Bayes

solution $S^*$ minimizes $E_\tau E_{\underline{\theta}} L(\underline{\theta}, S(\underline{X}_1, \ldots, \underline{X}_k))$ among all $S: \mathbb{R}^{nk} \to G$ by definition.

(Here and later on we tacitly assume that $L$ is intergrable properly.) The well

known standard method to find the Bayes rules $S^*$ is to minimize, after observing

$\underline{X}_i = \underline{x}_i$, $i = 1, \ldots, k$, the posterior expected loss $E(L(\underline{\theta}, s) \mid \underline{X}_i = \underline{x}_i$, $i = 1, \ldots, k)$

among all $s \in G$, which always is possible since $G$ is finite.

The first paper within this framework is due to Deely and Gupta (1968) and deals with the "linear" loss function $L(\underline{\theta},s) = \sum_{i \in s} \alpha_{s,i}(\theta_{[k]} - \theta_i)$, where the $\alpha$'s are non-negative and $\theta_{[1]} \leq \cdots \leq \theta_{[k]}$ denote the ordered values of $\theta_1, \ldots, \theta_k$. (This notation will be utilized analogously for all vectors used in the sequel.). The result is that, under an additional assumption on the $\alpha$'s, one can choose a Bayes rule which always selects only one population.

Because of this somewhat undesirable property Goel and Rubin (1977) choose the loss function $L(\underline{\theta},s) = c|s| + \theta_{[k]} - \max_{i \in s} \theta_i$ and study the behavior of the corresponding Bayes rule in great detail. (Because of the complexity of the problem it is necessary to derive approximate solutions.). Other papers dealing with other non-additive loss functions are due to Bickel and Yahav (1977), where the Bayesian aspect, however, is not of primary interest, and due to Chernoff and Yahav (1977), where two rules (one of these is Gupta's means procedure (c.f. (3.12))) are compared with the Bayes rule in a "normal model" on the basis of Monte Carlo results. As a result Gupta's maximum means procedure turns out to be "remarkably efficient" w.r.t. the Bayes rule.

Finally three papers dealing with additive (non-linear) loss functions of the type $L(\underline{\theta},s) = \sum_{i \in s} (c_2 - c_1 I_{\{\theta_{[k]}\}}(\theta_i))$ (where I denotes the indicator function) are due to Bratcher and Bhalla (1974) and Gupta and Hsu (1977,78). In the first the Bayes rule is derived and a binomial example is given. In the second similar Monte Carlo-studies are performed as in Chernoff and Yahav (1977), and again Gupta's maximum means procedures "do almost as well as the Bayes procedure". And in the last monotonicity of Bayes rules is the topic.

This paper serves two purposes: First (in Section 2) we show that the result of Deely and Gupta (1968) is not due to the linearity of the loss function but is due to the combined effect of the additivity of the loss

function and the non-negativity of its terms. Then we study the case of additive and especially linear loss functions, thereby filling a gap lying in between the paper by Deely and Gupta (1968) and the others mentioned above. Then (in Section 3) we show that in the normal case with symmetric normal priors Gupta's maximum means procedure turns out to be asymptotically Bayes w.r.t. a class of additive loss functions, whereas on the other hand Seal's procedure turns out to be Bayes w.r.t. an unrealistic additive loss function. Finally (in the appendix) we derive some bounds for $E(\max_{j=1,\ldots,k} (\mu_j + \rho V_j))$ (where $\underline{\mu} \in \mathbb{R}^k$, $\rho \in \mathbb{R}$ are fixed known and $\underline{V} \sim N(\underline{0}, I)$) to approximate the Bayes rules w.r.t. linear loss functions in cases where n is finite.

## 2. Additive and Linear Loss Functions.

As indicated above we first discuss the result of Deely and Gupta (1968). Let us consider the loss function $L(\underline{\theta}, s) = \sum_{i \in s} \alpha_{s,i} \ell_i(\underline{\theta})$ and assume for simplicity that $\alpha_{s,i} = \alpha(|s|)$, $i = 1, \ldots, k$, $s \in G$, holds. (In fact our following result remains valid if we replace this assumption by the corresponding one of Deely and Gupta (1968).). Then we can state the following slight generalization of the theorem in Deely and Gupta (1968):

Theorem 1. Let $m \, \alpha(m) \geq \alpha(1)$, $m = 1, \ldots, k$. If the $\ell_i$'s are non-negative then there exists a Bayes rule which always selects exactly one population.

Proof: Given $\underline{X} = \underline{x}$ the aposteriori risk of any procedure S is given by

$R(S|\underline{x}) = \alpha(|S(\underline{x})|) \sum_{i \in S(\underline{x})} A_i(\underline{x})$, where $A_i(\underline{x}) = E(\ell_i(\underline{\theta}), \underline{x}) \geq 0$, $i = 1, \ldots, k$.

Thus $R(S|\underline{x}) \geq \alpha(|S(\underline{x})|) |S(\underline{x})| \min_{i=1,\ldots,k} A_i(\underline{x})$

$$\geq \alpha(1) \min_{i=1,\ldots,k} A_i(\underline{x}) = \min_{i=1,\ldots,k} R(\{i\}|\underline{x}),$$

which completes the proof.

<u>Example 1.</u> That the converse statement does not hold true can be demonstrated by the following example: Let

$$L(\underline{\theta},s) = |s|^{-1} \sum_{i \in s} (\theta_{[k]} - \theta_i - \varepsilon)$$

$$= \theta_{[k]} - |s|^{-1} \sum_{i \in s} \theta_i - \varepsilon, \quad \varepsilon > 0.$$

Clearly $|S(\underline{x})|^{-1} \sum_{i \in S(\underline{x})} E(\theta_i|\underline{x}) \leq \max_{i=1,\ldots,k} E(\theta_i|\underline{x})$ holds, and even strict inequality occurs in many cases.

<u>Definition 1.</u> We call a loss function L <u>additive</u>, if

$$(2.1) \quad L(\underline{\theta},s) = \sum_{i \in s} \ell_i(\underline{\theta}), \quad \ell_i: \Omega^k \to \mathbb{R}, \quad i = 1,\ldots,k,$$

and <u>linear</u>, if

$$(2.2) \quad L(\underline{\theta},s) = c \sum_{i \in s} (\theta_{[k]} - \theta_i - \varepsilon), \quad \text{where } \varepsilon, c > 0, \text{ and } c \text{ clearly can be}$$

put equal to one.

In many situations the following assumptions assuring invariance under permutations and monotonicity of the losses seem to be quite natural:

(2.3) (a) $\ell_i(\underline{\theta}) = \ell(\{\theta_1,\ldots,\theta_k\}, \theta_i), \quad \underline{\theta} \in \Omega^k.$

    (b) $\ell_i(\underline{\theta}) \leq \ell_j(\underline{\theta})$ if $\underline{\theta} \in \Omega^k$ with $\theta_j \leq \theta_i$, $i,j \in \{1,\ldots,k\}.$

    (c) $\ell_i(\underline{\theta}) \geq \ell_i(\underline{\theta}')$ if $\underline{\theta}, \underline{\theta}' \in \Omega^k$ with $\theta_i \leq \theta_i'$ and

$$\theta_j \geq \theta_j', \quad j \neq i, \quad i = 1,\ldots,k.$$

Obviously these conditions are met by linear loss functions.

The following theorem, which to some extent can be viewed as being a special case of Lehmann's (1957) result, can be stated now without proof:

Theorem 2. For an additive loss function of type (2.1) the subsets $s(\underline{x}) \in G$ which a Bayes rule would select after observing $\underline{x} \in \mathbb{R}^{nk}$ satisfy the relation $\underline{S}^*(\underline{x}) \subseteq s(\underline{x}) \subseteq \tilde{S}^*(\underline{x}) \cup \{i\}$ for at least one $i \in M(\underline{x})$, where

$$(2.4) \quad \tilde{S}^*(\underline{x}) = \{j \mid E(\ell_j(\underline{\theta}) \mid \underline{x}) \leq 0\},$$

$$\underline{S}^*(\underline{x}) = \{j \mid E(\ell_j(\underline{\theta}) \mid \underline{x}) < 0\} \quad \text{and}$$

$$M(\underline{x}) = \{i \mid E(\ell_i(\underline{\theta}) \mid \underline{x}) = \min_{j=1,\ldots,k} E(\ell_j(\underline{\theta}) \mid \underline{x})\}.$$

Especially, for a linear loss function (2.4) reads as follows:

$$(2.5) \quad \tilde{S}^*(\underline{x}) = \{j \mid E(\Theta_j \mid \underline{x}) \geq E(\Theta_{[k]} \mid \underline{x}) - \varepsilon\},$$

$$\underline{S}^*(\underline{x}) = \{j \mid E(\Theta_j \mid \underline{x}) > E(\Theta_{[k]} \mid \underline{x}) - \varepsilon\} \quad \text{and}$$

$$M(\underline{x}) = \{i \mid E(\Theta_i \mid \underline{x}) = \max_{j=1,\ldots,k} E(\Theta_j \mid \underline{x})\}.$$

Remark 1. Another type of loss function which turns out to be equivalent to our type of additive loss function is the following:

$$(2.6) \quad \hat{L}(\underline{\theta}, s) = \sum_{i \in s} \ell_i^+(\underline{\theta}) + \sum_{i \notin s} \ell_i^-(\underline{\theta}),$$

where $\ell_i^+, \ell_i^- : \Omega^k \to [0, \infty)$, $i = 1, \ldots, k$, can be viewed as being losses for errors of the first and second kind in analogy to testing theory. Since such a loss function can be rewritten as

$$(2.7) \quad \hat{L}(\underline{\theta}, s) = \sum_{i=1}^{k} \ell_i^-(\underline{\theta}) + \sum_{i \in s} (\ell_i^+(\underline{\theta}) - \ell_i^-(\underline{\theta})),$$

we arrive at the same set of Bayes solutions if we drop the first sum on the r.h.s. of (2.7). (Hereby in fact the overall risk is changed only by an additive constant. Note that the change of the conditional risk, given $X = \underline{x}$, has no influence on the determination of the Bayes rules.).

Conversely if we start with an additive loss function of type (2.1) we can switch over to a loss function of type (2.6) analogously by choosing

$\ell_i^+$ ($\ell_i^-$) to be the usual positive (negative) part of $\ell_i$, $i = 1,\dots,k$.

Thus Theorem 1 - applied to (2.6) - can be interpreted as follows: if we always have to pay for every population selected, then we take as few as possible: namely always one.

Example 2. Bratcher and Bhalla took a loss function $\hat{L}$ of type (2.6) with

$$\ell_i^+(\underline{\theta}) = \hat{c}_1(1-I_{\{\theta_{[\kappa]}\}}(\ell_i)), \quad \ell_i^-(\underline{\theta}) = \hat{c}_2 I_{\{\theta_{[\kappa]}\}}(\theta_i), \quad \hat{c}_1,\hat{c}_2 > 0, \quad i = 1,\dots,k.$$

On the other hand Gupta and Hsu (1977,78) took the loss function $L(\underline{\theta},s) =$

$c_1(1 - I_{\{\theta_i | i \in s\}}(\theta_{[k]})) + c_2|s|$, $c_1$, $c_2 > 0$.

If we assume that the posterior distribution of $\Theta$ assures that for every $\underline{x} \in \mathbb{R}^{nk}$ all the $\Theta_i$'s are distinct with probability one, then $L(\underline{\theta},s)$ can be replaced by

$$L(\underline{\theta},s) = c_1 + \sum_{i \in s} (c_2 - c_1 I_{\{\theta_{[k]}\}}(\theta_i)).$$

If we take into account that we can multiply a loss function with a positive constant and moreoever can add any further constant to it without changing a given Bayes problem, then it is easy to see that the loss function of Bratcher and Bhalla (1974) and that of Gupta and Hsu (1978) are equivalent in the sense of Remark 1 and, especially, both are additive.

It seems to be worth mentioning that in the normal model (cf. (3.1)) the converse of Theorem 1 also holds true.

If the $\ell_i$'s mentioned there may also assume negative values (i.e. if $c_2 < c_1$ here), then there exists no Bayes rule which always selects exactly one population. Since aposteriori the $\Theta_i$'s now are jointly (non-degenerate) normally distributed, this follows from the fact that the Bayes rule S*, for all $\underline{x} \in \mathbb{R}^{nk}$ except possibly a null set, turns out to be

(2.8)   $i \in S^*(\underline{x})$ iff $P\{\Theta_i = \Theta_{[k]}|\underline{x}\} = \max_{j=1,\ldots,k} P\{\Theta_j = \Theta_{[k]}|\underline{x}\}$

and/or $P\{\Theta_i = \Theta_{[k]}|\underline{x}\} > c_2/c_1$.

Finally let us consider some other properties of the Bayes rules for additive loss functions. For this purpose let $Z_i = H(\underline{X}_i)$ be a one-dimensional sufficient statistic for $\theta_i$ with distribution $Q_{\theta_i}$, say, $i = 1,\ldots,k$. The rules S of course depend now on $Z_1,\ldots,Z_k$ and are defined on $\mathbb{R}^k$.

<u>Definition 2</u>. We call a rule $S: \mathbb{R}^k \to G$ <u>ordered</u> if for every

$\underline{z} \in \mathbb{R}^k$ $i \in S(\underline{z})$ and $z_i < z_j$ implies $j \in S(\underline{z})$.

We call it <u>monotone</u> if for every $i \in \{1,\ldots,k\}$ and $\underline{z}, \underline{z}' \in \mathbb{R}^k$ with

$z_i \leq z_i'$ and $z_j \geq z_j'$, $j \neq i$, $i \in S(\underline{z})$ implies $i \in S(\underline{z}')$.

Since we only need to consider non-randomized rules in this Bayesian framework (cf. Goel and Rubin (1977)), Definition 2 for monotonicity and Definition 8 in Gupta and Hsu (1977) coincide. Now we can state:

<u>Theorem 3</u>.   <u>Let</u> $\{Q_\theta\}_{\theta \in \Omega}$ <u>have</u> <u>densities</u> $\{f_\theta\}_{\theta \in \Omega}$ <u>w.r.t. the Lebesgue measure</u> <u>on</u> $(\mathbb{R},\mathcal{B})$ <u>with monotone non-decreasing likelihood ratios in</u> $\theta$, <u>and let</u> L <u>be any additive loss function</u>. Then the following statements hold true:

(i)   <u>If</u> L <u>satisfies</u> (2.3) (a) <u>and</u> (b) <u>and if the prior distribution of</u> $\Theta$ <u>is</u> <u>symmetric on</u> $\Omega^k$, <u>then every Bayes rule can be assumed to be ordered</u>.

(ii)   <u>If</u> L <u>satisfies</u> (2.3) (c) <u>and if</u>

$\Theta_1,\ldots,\Theta_k$ <u>are apriori independently distributed, then every Bayes rule</u> <u>can be assumed to be monotone</u>.

<u>Proof</u>: The first part follows from Goel and Rubin (1977) (cf. Lemma 1 and Remark 1 cited there), since our loss function satisfies the conditions stated there.

The second part is a generalization of the Theorem 3 of Gupta and Hsu (1977). It is easy to see that their proof works with every loss function whose components $\ell_1, \ldots, \ell_k$ satisfy the conditions stated in (II).

Note that additive loss functions of type (2.1) with $\ell_i(\underline{\theta}) = \ell(\theta_{[k]} - \theta_i)$, $i = 1, \ldots, k$, for non-decreasing $\ell$ and thus in particular linear loss functions satisfy all the conditions stated in the theorem above. Thus we can state:

Corollary 1. <u>Let</u> $\{Q_\theta\}_{\theta \in \Omega}$ <u>be given as in Theorem 3 and let</u>

$$L(\underline{\theta}, s) = \sum_{i \in s} \ell(\theta_{[k]} - \theta_i) \quad \underline{\text{with non-decreasing}}$$

$\ell: \mathbb{R} \to \mathbb{R}$. <u>Then, if</u> $\Theta_1, \ldots, \Theta_k$ <u>are apriori independently</u> <u>identically distributed, every Bayes rule can be assumed to</u> <u>be ordered and monotone.</u>

Remark 2. In their Lemma 2 Goel and Rubin (1977) give an aid with which one can simplify the computation of the Bayes procedure if one adopts their loss function and assumes a symmetric prior distribution:

If $\underline{Z} = \underline{z} \in \mathbb{R}^k$ is observed and $\underline{z}$ is, without loss of generality, ordered in such a way that $z_1 \leq \ldots \leq z_k$ holds, then the terms $r_i(\underline{z}) = E(L(\underline{\Theta}, \{k, k-1, \ldots, k-i+1\}) | \underline{z})$, $i = 1, \ldots, k$, have the property that $r_{j+1}(\underline{z}) - r_j(\underline{z})$ is non-decreasing in $j = 1, 2, \ldots, k-1$.

For additive loss functions this property reduces just to

$$(2.9) \quad E(\ell_1(\underline{\theta}) | \underline{z}) \geq \ldots \geq E(\ell_k(\underline{\theta}) | \underline{z}), \quad z_1 \leq \ldots \leq z_k.$$

Though it is difficult to find general sufficient conditions for (2.9) to hold true, we can at least state the following:

Theorem 4. <u>Let</u> $\{Q_\theta\}_{\theta \in \Omega}$ <u>and L be given as in Corollary 1. Then each of the</u> <u>two conditions stated below is sufficient for (2.9) to hold true:</u>

(A.1)   $\Theta_1, \ldots, \Theta_k$ _are apriori independently identically distributed._

(A.2)   L _is linear and the prior distribution_ $\tau$ _is symmetric on_ $\Omega^k$.

_Proof:_ Under (A.1) note that, given $\underline{Z} = \underline{z}$ with $z_1 \leq \ldots \leq z_k$, the $\Theta_i$'s are independent and stochastically ordered in the same order as the $z_i$'s. With standard analysis one can show that this implies that even the dependent variables $\Theta_1 - \Theta_{[k]}, \ldots, \Theta_k - \Theta_{[k]}$ are stochastically ordered in the same direction. Thus by the monotonicity of $\ell$ the assertion follows.

Under (A.2) for every $\underline{z} \in \mathbb{R}^k$ with $z_1 \leq \ldots \leq z_k$ and $i \in \{1, \ldots, k-1\}$ we have

$$E(\ell_i(\underline{\Theta}) | \underline{z}) - E(\ell_{i+1}(\underline{\Theta}) | \underline{z}) = E(\Theta_{i+1} - \Theta_i | \underline{z}) = b \int_{\Omega^k} (\Theta_{i+1} - \Theta_i) \prod_{r=1}^{k} f_{\Theta_r}(z_r) d\tau(\underline{\Theta})$$

$$= b \int_{\{\Theta_i < \Theta_{i+1}\}} (\Theta_{i+1} - \Theta_i) [f_{\Theta_i}(z_i) f_{\Theta_{i+1}}(z_{i+1}) - f_{\Theta_i}(z_{i+1}) f_{\Theta_{i+1}}(z_i)] \prod_{r \neq i, i+1} f_{\Theta_r}(z_r) d\tau(\underline{\Theta}),$$

where b is a normalizing factor and the last identity follows by the method which Goel and Rubin used to prove their Lemma 2. Thus by the M.L.R. property of $\{f_\theta\}_{\theta \in \Omega}$ the proof is completed.

## 3. The Normal Model.

In this section we assume that apriori $\underline{X}_1, \ldots, \underline{X}_k$ are independent samples of common size n from k normal populations with unknown means $\theta_1, \ldots, \theta_k$ and a common known variance $\sigma^2 > 0$. By sufficiency we can reduce our set of data to $\underline{X} = (\bar{X}_1, \ldots, \bar{X}_k)$, where $\bar{X}_1, \ldots, \bar{X}_k$ are the corresponding sample means. As to the $\Theta$'s we assume an exchangeable normal prior. More precisely our "normal model" is as follows:

(3.1)   (a)   $\underline{X} | \underline{\Theta} = \underline{\theta} \sim N(\underline{\theta}, q\ I)$   and

      (b)   $\underline{\Theta} \sim N(m\ \underline{1}, r\ I + t\ U)$, where

      $q = \sigma^2/n, \ m \in \mathbb{R}, \ r > 0, \ t > -r/k, \ \underline{1} = (1, \ldots, 1)'$,

      $U = \underline{1}\ \underline{1}'$ and 1 denotes the k×k identity matrix.

Note that $r > 0$ together with $t > -r/k$ is necessary and sufficient for

r I + t U to be positive definite. This model was chosen by Chernoff and Yahav (1977) (with $t > 0$) and by Gupta and Hsu (1978).

By (3.1) it is easy to see that we have

(3.2) (a) $\underline{\theta} | \underline{X} = \underline{x} \sim N(\underline{\mu}, a\,I + b\,U)$, where

$$\underline{\mu} = r(q+r)^{-1}\,\underline{x} + q\,t((q+r)(q+r+kt))^{-1}\,U\underline{x} + q(q+r+kt)^{-1}\,m\,\underline{1},$$

$$a = r\,q\,(q+r)^{-1} \quad \text{and} \quad b = q^2 t((q+r)(q+r+kt))^{-1}, \text{ and}$$

(b) $\underline{X} \sim N(\underline{\tilde{\mu}}, \tilde{a}\,I + \tilde{b}\,U)$, where

$$\underline{\tilde{\mu}} = m\,\underline{1}, \quad \tilde{a} = q+r \quad \text{and} \quad \tilde{b} = t.$$

Let us include non-additive loss functions $L(\underline{\theta},s)$ into our next considerations. Then one intuitively feels that one can put m and t equal to zero without changing the problem, if L is translation-invariant, i.e. if $L(\underline{\theta},s) = L(\underline{\theta} + \eta\,\underline{1}, s)$ holds for all $\underline{\theta} \in \Omega^k$, $s \in G$ and $\eta \in \mathbb{R}$. Moreover one should expect that the Bayes rules then are translation-invariant, too: $S^*(\underline{x}) = S^*(\underline{x} + \eta\,\underline{1})$ for all $\eta \in \mathbb{R}$ and all $\underline{x} \in \mathbb{R}^k$ except possibly a null set.

Ideas of this kind primarily are due to Chernoff and Yahav (1977) and also to Gupta and Hsu (1977). Nevertheless since the formulation and proof of a general theorem in this direction is missing up to now we feel that it is justified to do this in our present paper.

Now if any random vector $\underline{Y}$ is distributed according to some $N(\underline{\mu}, a\,I + b\,U)$ with $a > 0$ and $b > -a/k$, then at once one has in mind the following ("conditional i.i.d.") representation:

(3.3) $\underline{Y} = a^{1/2}\,\underline{V} + b^{1/2}\,W\,\underline{1} + \underline{\mu}$ where

$\underline{V} \sim N(\underline{0},I)$ and $W \sim N(0,1)$ are independent.

But this holds only in cases where $b \geq 0$ (i.e. where the correlations of the Y's are non-negative), and therefore does not help us in our more general setup. But, fortunately, there exists another representation, which always can be used for our purposes:

<u>Lemma 1</u>. <u>Let</u> $\underline{\mu} \in \mathbb{R}^k$, $a > 0$ <u>and</u> $b > -a/k$. <u>If</u>

$$(3.4) \qquad \underline{Y} = a^{1/2} \underline{V} + k^{-1}((a+k\,b)^{1/2} - a^{1/2}) U \underline{V} + \underline{\mu},$$

$$\underline{\text{with}}\ \underline{V} \sim N(\underline{0},\ I),$$

$$\underline{\text{then}}\ \underline{Y} \sim N(\underline{\mu},\ a\,I + b\,U).$$

The proof is standard and therefore omitted. Besides we remark that in (3.4) $(a + k\,b)^{1/2}$ can be replaced by $-(a + k\,b)^{1/2}$.

<u>Theorem 5</u>. <u>Under the normal model (3.1) for every loss function</u> $L(\theta,s)$ <u>which is translation-invariant the following three assertions hold true:</u>

(i) <u>For every rule</u> $S: \mathbb{R}^k \to G$ <u>and every</u> $\underline{x} \in \mathbb{R}^k$ <u>the posterior risk</u> $E(L(\theta,S(\underline{x})) | \underline{X} = \underline{x})$ <u>does not depend on m</u> <u>and</u> t.

(ii) <u>Every Bayes rule can be assumed to be translation-invariant.</u>

(iii) <u>For every translation-invariant rule</u> S <u>the overall risk</u> $E(L(\theta,S(\underline{X})))$ <u>does not depend on m</u> <u>and</u> t.

<u>Proof</u>: Under the normal model (3.1) let $L(\underline{\theta},s)$ be translation-invariant.

(i) Given $\underline{X} = \underline{x}$, $\underline{\theta}$ is distributed according to $N(\underline{\mu},\ a\,I + b\,U)$ where $\underline{\mu}$, $a$ and $b$ are given by (3.2) (a). In this situation we choose for $\underline{\theta}$ the representation given by Lemma 1.

Let S be any rule, $S(\underline{x}) = s$, say, and $\rho \in \mathbb{R}$. Then since $U\,\underline{y} = (y_1 + \ldots + y_k)\,\underline{1}$ for every $\underline{y} \in \mathbb{R}^k$, we get by the translation-invariance of $L(\underline{\theta},s)$

$$(3.5) \quad E(L(\underset{\sim}{\theta},s)|\underset{\sim}{X} = \underset{\sim}{x} + \rho \underset{\sim}{1}) = E(L(r(q+r)^{-1}\underset{\sim}{x} + (rq(q+r)^{-1})^{1/2} \underset{\sim}{v},s)).$$

Since the r.h.s. does not depend on m, t and $\rho$, part (I) is proved by putting $\rho = 0$.

(ii) If, given $\underset{\sim}{X} = \underset{\sim}{x}$, one $s^* \in G$ minimizes the l.h.s. of (3.5) for one $\rho \in \mathbb{R}$, then it does it for all $\rho \in \mathbb{R}$ simultaneously. Thus every Bayes rule is translation-invariant, if one neglects possible pathological choices in cases where several solutions appear.

(iii) Let S be a translation-invariant rule and let for $\underset{\sim}{x} \in \mathbb{R}^k$

$$R(S|\underset{\sim}{x}) = E(L(\underset{\sim}{\theta},S(\underset{\sim}{x}))|\underset{\sim}{X} = \underset{\sim}{x})$$

denote the conditional risk of S - given $\underset{\sim}{X} = \underset{\sim}{x}$. By (3.5) we see that $R(S|\underset{\sim}{x})$ does not depend on m and t and moreover is translation-invariant in $\underset{\sim}{x}$.

Now, marginally, $\underset{\sim}{X}$ is distributed according to $N(\underset{\sim}{\tilde{\mu}},\tilde{a} I + \tilde{b} U)$ where $\underset{\sim}{\tilde{\mu}}$, $\tilde{a}$ and $\tilde{b}$ are given by (3.2) (b). If we choose now for $\underset{\sim}{X}$ the representation given by Lemma 1 then we see that the overall risk $E(L(\underset{\sim}{\theta}, S(\underset{\sim}{X}))) = E(R(S|\underset{\sim}{X}))$ does not depend on m and t. Thus the proof of the theorem is completed.

Remark 3. In the present (Bayesian) framework Theorem 5 clearly fits all our needs. But it should be pointed out that the following simple fact can be viewed as being the basis of this theorem:
"Let $\underset{\sim}{Y} \sim N(\underset{\sim}{\mu} + \rho \underset{\sim}{1}, a I + b U)$ with $\underset{\sim}{\mu} \in \mathbb{R}^k$, $\rho \in \mathbb{R}$, $a > 0$ and $b > -a/k$. Then there exists a random vector $\underset{\sim}{Z}$ with $\underset{\sim}{Z} \sim N(\underset{\sim}{\mu}, a I)$ such that $h(\underset{\sim}{Y}) = h(\underset{\sim}{Z})$ everywhere for every translation-invariant $h: \mathbb{R}^k \to \mathbb{R}^k$".

For the remainder of this paper we restrict our considerations to additive loss functions with translation-invariant $\ell_i(\underset{\sim}{\theta})$, $i = 1,...,k$. Since by Theorem 5 we can put $m = t = 0$ without loss of generality, our model can considerably be simplified to

(3.6)  $\underline{X}|\underline{\theta} = \underline{\theta} \sim N(\underline{\theta}, q\ I)$,  $\underline{\theta} \sim N(\underline{0},\ r\ I)$,

which in turn implies

(3.7)  $\underline{\theta}|\underline{X} = \underline{x} \sim N(r(q+r)^{-1}\underline{x},\ r\ q(q+r)^{-1}\ I)$,  $\underline{X} \sim N(\underline{0},(q+r)I)$.

Thus, given $\underline{X} = \underline{x}$, $\underline{\theta}$ has the representation:

(3.8)  $\underline{\theta} = r(q+r)^{-1}\underline{x} + (r\ q(q+r)^{-1})^{1/2}\ \underline{V}$ with $\underline{V} \sim N(\underline{0}, I)$.

The next result is to some extent sharper than Theorems 3 and 4:

Theorem 6.  If under the normal model (3.1)

$$L(\underline{\theta}, s) = \sum_{i\ \in\ s} \ell\ (\theta_{[k]} - \theta_i)$$ with nondecreasing $\ell$, then (2.9) holds

and every Bayes rule can be assumed to be ordered and monotone.

Proof:  By Theorem 2 every Bayes rule selects according to small values of

$E(\ell(\theta_{[k]} - \theta_i)|\underline{x})$.  Since $\ell$ is non-decreasing, we have for every $\underline{\theta} \in \Omega^k$

$\ell(\theta_{[k]} - \theta_i) = \max_{j=1,\ldots,k} \ell(\theta_j - \theta_i)$, $i = 1,\ldots,k$.  But by (3.8) we get for

every $\underline{x} \in \mathbb{R}^k$

(3.9)  $E(\ell(\theta_{[k]} - \theta_i)|\underline{x}) = E(\max_{j=1,\ldots,k} \ell[r(q+r)^{-1}(\bar{x}_j - \bar{x}_i) + (r\ q(q+r)^{-1})^{1/2}$

$$(V_j - V_i)]),$$

which clearly implies the desired result.

Example 3.  Under the normal model (3.1) consider the additive loss function

(3.10)  $L(\underline{\theta}, s) = \sum_{i\ \in\ s} [k^{-1} \sum_{j=1}^{k} \theta_j - \theta_i - \epsilon]$, $\epsilon > 0$.

Since it is translation-invariant, by Theorem 4 we can assume that (3.8) holds,

and the unique Bayes rule S* turns out to be

(3.11)  $i \in S^*(\underline{X})$ iff $\bar{X}_i \geq k^{-1} \sum_{j=1}^{k} \bar{X}_j - r^{-1}(q+r)\epsilon$,

which is the well known procedure of Seal (1957).

There is a long story about the question how good different well established procedures like this perform in certain circumstances. One result is that Seal's procedure is not safe to use and especially inferior to Gupta's maximum means procedures (cf. Definition 3) with respect to many aspects. This is shown for example in Seal (1957), Deely and Gupta (1968) and Gupta and Miescke (1978). One new argument in this direction now seems to be to us that Seal's procedure under the normal model turns out to be a Bayes solution w.r.t. a very unrealistic loss function.

Perhaps the best known subset selection procedure is due to Gupta (1956, 65):

Definition 3. Gupta's maximum means procedure is given by

$$(3.12) \quad i \in S_d(\underline{X}) \text{ iff } \bar{X}_i \geq \bar{X}_{[k]} - q^{1/2} d, \, d > 0.$$

In the classical (non-Bayesian) approach due to Gupta one has to choose an $P^* > k^{-1}$ and then $d(k,P^*)$ is given by the requirement that $S_d$ should contain the best population with probability at least $P^*$ for every fixed $\underline{\theta} \in \Omega^k$. (Conversely if $d$ is predetermined then $P^*(k,d)$ is fixed.)

It is conjectured and partially proved by many authors that $S_d$ performs well or even is optimal in many situations. But up to now it was not possible to find $S_d$ to be close to a Bayes rule in any given model, except perhaps in the Monte Carlo-studies of Chernoff and Yahav (1977) and Gupta and Hsu (1977). In this spirit our following results seem to be interesting.

Theorem 7. If under the normal model (3.1)

$$L(\underline{\theta},s) = \sum_{i \in s} \ell(\theta_{[k]} - \theta_i - \epsilon), \text{ where } \epsilon > 0 \text{ is fixed}, \ell \text{ is non-}$$
decreasing, continuous, bounded and satisfies $\ell(\rho) = 0$ if and only if $\rho = 0$, then the following procedure S is the limit of Bayes rules for large n:

(3.13)     $i \in S(\underline{X})$ __iff__ $\bar{X}_i \geq \bar{X}_{[k]} - \epsilon$.

__Proof:__ Given $\underline{X} = \underline{x}$, by Theorem 2 and 5 and especially (3.8) the following rules are Bayes rules for $n = 1,2,\ldots$:

$$i \in S_n^*(\underline{x}) \text{ iff } E(\ell[\max_{j=1,\ldots,k} \{r(q+r)^{-1}(\bar{x}_j - \bar{x}_i)$$

$$+ (r\,q\,(q+r)^{-1})^{1/2}\,(V_j - V_i)\} - \epsilon]) \leq 0.$$

Since $\ell$ is bounded and continuous, for large $n$ (by Lebesgue's dominated convergence theorem) the expectation converges to $\ell[\max_{j=1,\ldots,k} \{\bar{x}_j - \bar{x}_i\} - \epsilon]$.

Thus in view of the additional assumptions which we imposed upon $\ell$ the theorem is proved.

Though we have seen that the limit of Bayes rules is very similar to rules of type (3.12), this does not completely satisfy our requests, since procedures of type (3.13) have a $P^* = P^*(k, q^{-1/2} \epsilon)$ which tends to one for large $n$. Thus for a moment we alternatively take another loss function as given in the theorem below:

__Theorem 8.__  __If under the normal model (3.1)__

$L(\underline{\theta}, s) = \sum_{i \in s} (\theta_{[k]} - \theta_i - q^{1/2} d)$ __with a fixed__ $d > 0$, __then Gupta's maximum means__ __procedure__ $S_d$ __is the limit of the (unique) Bayes rules as__ $n$ __tends to infinity.__

__Proof:__ Given $\underline{X} = \underline{x}$, by Theorems 2 and 5 and especially (3.8) the unique Bayes rules for $n = 1,2,\ldots$ turn out to be

(3.14)        $i \in S_n^*(\underline{x})$   iff

$$\bar{x}_i \geq E(\max_{j=1,\ldots,k} \{\bar{x}_j + (r^{-1}q(q+r))^{1/2} V_j\}) - r^{-1}(q+r)q^{1/2} d,$$

if the $\bar{x}_j$'s are distinct. By (A.9) we have

(3.15) $\quad \bar{x}_{[k]} \leq E(\max_{j=1,..,k} \{\bar{x}_j + (r^{-1}q(q+r))^{1/2} V_j\}) \leq \bar{x}_{[k]} + \tilde{A}_2(\underline{x})$, where

$$\tilde{A}_2(\underline{x}) = (2r^{-1}q(q+r))^{1/2} \sum_{j=1}^{k-1} T((2r^{-1}q(q+r))^{-1/2}(\bar{x}_{[j]} - \bar{x}_{[k]}))$$

and the function T is given by (A.1).

Since by (3.7) $\underline{X}$, unconditionally, has non-degenerate normal distribution, we can assume that all the $\bar{x}_j$'s are distinct. But then by (A.2) we have $\tilde{A}_2(\underline{x}) = o(q^{1/2})$ and the proof is completed by noting that r remains fixed and $q = \delta^2/n$ tends to zero if n tends to infinity.

At least we shall study the case of linear loss functions in more detail. Here we have to distinguish between two possibilities:

(3.16) $\quad L_1(\underline{\theta},s) = \sum_{i \in s} (\theta_{[k]} - \theta_i - \varepsilon)$, with fixed $\varepsilon > 0$ and

(3.17) $\quad L_2(\underline{\theta},s) = \sum_{i \in s} (\theta_{[k]} - \theta_i - q^{1/2} d)$, with fixed $d > 0$.

By Theorems 7 and 8 we know that asymptotically $L_2$ leads to exactly one population (the best in every situation where $\underline{\theta} = \theta$ is fixed) whereas $L_1$ leads to a screening procedure even in the limit. On the other hand for every finite n both $L_1$ and $L_2$ give us screening procedures.

For the case of finite n we give now some approximations to the Bayes rule which apply to $L_1$ and $L_2$ simultaneously. These approximations are very easy to handle since they do not involve any integral and can in fact be evaluated with the help of the function T alone. For convenience we formulate our results in terms of $L_1$.

The results for $L_2$ follow easily since one has to replace $\varepsilon$ by $q^{1/2}d$ only.

Corollary 2. If under the normal model (3.1) $L_1(\underline{\theta},s)$ is given by (3.16) and $\varepsilon > 0$ and n are fixed, then with probability one we have for $\alpha = 1,2$ and $\beta = 1,\ldots,4$

(3.18)   $S_\alpha(\underline{X}) \subseteq S^*(\underline{X}) \subseteq \tilde{S}_\beta(\underline{X})$,

    <u>where for</u> $\underline{x} \in \mathbb{R}^k$ $i \in S_\alpha(\underline{x})$ $(\underline{or}\ \tilde{S}_\beta(\underline{x}))$ <u>iff</u>

$$\bar{x}_i \geq \bar{x}_{[k]} - r^{-1}(q+r)\varepsilon + \tilde{A}_\alpha(\underline{x})\ (\underline{or}\ A_\beta(\underline{x})),$$

<u>and</u> $\tilde{A}_\alpha$, $A_\beta$ <u>are given by</u> (A.8-11) <u>and</u> (A.17-18).

<u>Proof:</u> The assertion follows from (3.14), where $q^{1/2} d$ is to be replaced by $\varepsilon$, and by the fact that for all $\alpha, \beta$ and $\underline{x} \in \mathbb{R}^k$ we have

$$\bar{x}_{[k]} + A_\beta(\underline{x}) \leq E(\max_{j=1,\ldots,k}(\bar{x}_j + (r^{-1}q(q+r))^{1/2} V_j)) \leq \bar{x}_{[k]} + \tilde{A}_\alpha(\underline{x}),$$

which is proved in the appendix.

    Let us explicitly point out the following special case of (3.18) if we take $\tilde{A}_1$ and $A_1$: Then with probability one we have

(3.19)   $S_{d_1}(\underline{X}) \subseteq S^*(\underline{X}) \subseteq S_{d_2}(\underline{X})$,

where $d_2 = q^{-1/2} r^{-1}(q+r)\varepsilon$, $d_1 = q^{-1/2} r^{-1}(q+r)\varepsilon - (r^{-1}(q+r))^{1/2} a_k$,

and $a_k$ is given by (A.2).

    Here we can expect that for moderate $n$ in many cases of $\underline{X} = \underline{x}$ $S_{d_1}$ and $S_{d_2}$ coincide, so that the experimenter finds the Bayes rule with the help of two means procedures, which are very easy to compute. This idea of course analogously applies to the other approximations, which, however, are no longer means procedures of type (3.12). Finally it should be pointed out that it is always possible to use the exact Bayes procedure (3.14) (with $q^{1/2} d$ possibly replaced by $\varepsilon$) if one is willing to evaluate either (A.4) or (A.5) with the help of a computer program.

Appendix: Expectation of the Largest of k Independently Distributed Normal

Random Variables.

We derive now lower and upper bounds for

$$E(\max_{j=1,\ldots,k} (\mu_j + \rho V_j)) = E_k(\underline{\mu}, \rho), \text{ say, where}$$

$\mu_j = \bar{x}_j$, $j = 1, \ldots, k$ and $\rho^2 = r^{-1} q(q+r)$ are fixed known and

$$\underline{V} = (V_1, \ldots, V_k)' \sim N(\underline{0}, I).$$

Remark 4. Let us mention briefly that

$$r(q+r)^{-1} E(\max_{j=1,\ldots,k} (\mu_j + \rho V_j))$$

turns out to be Bayes estimate for the largest mean (i.e. $\theta_{[k]}$) under model

(3.6) w.r.t. squared error loss. Thus our bounds (to be derived in the sequel)

multiplied with $r(q+r)^{-1}$ can be utilized also as bounds for these Bayes estimates.

Let $\varphi$ and $\Phi$ denote the density and distribution function of the one-

dimensional standard normal distribution, and let T be the following auxiliary

function which previously was used also by Goel and Rubin (1977):

$$(A.1) \quad T(\xi) = \int_{-\infty}^{\xi} \Phi(\eta) d\eta = \varphi(\xi) + \xi\Phi(\xi), \quad \xi \in \mathbb{R}.$$

T is strictly increasing, strictly convex and satisfies

$$(A.2) \quad \lim_{\xi \to -\infty} T(\xi) = 0, \quad T(0) = (2\pi)^{-1/2}, \quad \lim_{\xi \to \infty} (T(\xi) - \xi) = 0.$$

Moreover let

$$(A.3) \quad a_k = E(\max_{j=1,\ldots,k} V_j), \quad k = 1, 2, \ldots .$$

For further study and especially tables of the $a_k$'s see David (1970).

Lemma 2.

(A.4)  $E_k(\underline{\mu},\rho) = \mu_{[k]} + \sum_{j=1}^{k-1} \int_R \prod_{i=1}^{j} \Phi(\rho^{-1}(\xi-\mu_{[k-i+1]}))[1-\Phi(\rho^{-1}(\xi-\mu_{[k-j]}))]d\xi$

(A.5)  $E_k(\underline{\mu},\rho) = \mu_{[k]} +$

$$\int_R \Phi(\rho^{-1}(\xi-\mu_{[k]}))[1- \prod_{i=1}^{k-1} \Phi(\rho^{-1}(\xi-\mu_{[i]}))]d\xi.$$

Proof:  By Chernoff and Yahav (1977) we have the following recursive relations for $j = 1,\ldots,k-1$:

$$E(\max_{i \geq k-j} (\mu_i + \rho V_i)) - E(\max_{i \geq k-j+1} (\mu_i + \rho V_i))$$

$$= \int_{IR} \prod_{i=1}^{j} \Phi(\rho^{-1}(\xi-\mu_{[k-i+1]}))[1-\Phi(\rho^{-1}(\xi-\mu_{[k-j]}))]d\xi,$$

which clearly imply (A.4). And (A.5) follows by the telescopic property of the sum in (A.4).

For $k = 2$ this reduces to

Lemma 3.

(A.6)  $E_2(\underline{\mu},\rho) = \mu_{[2]} + 2^{1/2}\rho \ T(2^{-1/2} \rho^{-1}(\mu_{[1]}-\mu_{[2]}))$.

Proof:

$$E_2(\underline{\mu},\rho) = E(\mu_{[2]} + \rho V_2 + \max(0,\mu_{[1]} + \rho V_1 - \mu_{[2]} - \rho V_2))$$

$$= E(\mu_{[2]} + \rho V_2 + \max(0,\mu_{[1]} - \mu_{[2]} + 2^{1/2} \rho V_1))$$

$$= \mu_{[2]} + E((\mu_{[1]} - \mu_{[2]} + 2^{1/2} \rho V_1)^+).$$

Applying Lemma 7 of Goel and Rubin (1977) the proof is completed. Besides we remark that (A.6) in another context (cf. Remark 4) was also derived by Blumenthal and Cohen (1968).

For later applications let us rewrite Lemma 3 in the following way:

$$(A.7) \qquad \int_{\mathbb{R}} \Phi(\rho^{-1}(\xi+\Delta))\Phi(-\xi^{-1}\xi)d\xi = 2^{1/2}\rho \ T(2^{-1/2}\rho^{-1}\Delta) \text{ for all } \Delta \in \mathbb{R}.$$

Then we state:

Lemma 4. $E_k(\underline{\mu},\rho) \leq \mu_{[k]} + \tilde{A}_\alpha$, $\alpha = 1,2$, where

$$(A.8) \qquad \tilde{A}_1 = \rho \ a_k, \quad \text{and}$$

$$(A.9) \qquad \tilde{A}_2 = 2^{1/2}\rho \sum_{j=1}^{k-1} T(2^{-1/2}\rho^{-1}(\mu_{[j]} - \mu_{[k]})).$$

Proof: (A.8) is immediate. In view of (A.4) we have

$$E_k(\underline{\mu},\rho) \leq \mu_{[k]} + \sum_{j=1}^{k-1} \int_{\mathbb{R}} \Phi(\rho^{-1}(\xi-\mu_{[k]}))\Phi(-\rho^{-1}(\xi-\mu_{[k-j]}))d\xi.$$

Application of (A.7) and reordering of indices thus leads to (A.9).

Lemma 5. $E_k(\underline{\mu},\rho) \geq \mu_{[k]} + \underset{\sim}{A}_\beta$, $\beta = 1,2$, where

$$(A.10) \qquad \underset{\sim}{A}_1 = 0, \quad \text{and}$$

$$(A.11) \qquad \underset{\sim}{A}_2 = 2^{1/2}\rho \ T(2^{-1/2}\rho^{-1}(\mu_{[k-1]} - \mu_{[k]})).$$

Proof: (A.10) is immediate. In view of (A.5) we have

$$E_k(\underline{\mu},\rho) \geq \mu_{[k]} + \int_{\mathbb{R}} \Phi(\rho^{-1}(\xi-\mu_{[k]}))\Phi(-\rho^{-1}(\xi-\mu_{[k-1]}))d\xi,$$

and thus application of (A.7) leads to (A.11).

By the convexity of T we get two further bounds with the help of Jensen's inequality. At first we state

Lemma 6.

$$(A.12) \qquad E_k(\underline{\mu},\rho) = \mu_{[k]} + \rho E(T(\rho^{-1}[\underset{j=1,\ldots,k-1}{\max}(\mu_j+\rho \ V_j)-\mu_{[k]}]))$$

$$(A.13) \quad E_k(\underline{\mu},\delta) \geq \mu_{[k]} + \rho \, T(\rho^{-1}[E(\max_{j=1,\ldots,k-1}(\mu_j + \rho \, V_j)) - \mu_{[k]}]).$$

Proof: Integrating (A.5) by parts leads to

$$(A.14) \quad E_k(\underline{\mu},\rho) = \mu_{[k]} + \rho \int_{\mathbb{R}} T(\rho^{-1}(\xi - \mu_{[k]})) d\{ \prod_{i=1}^{k-1} \Phi(\rho^{-1}(\xi - \mu_{[i]})) \},$$

which equals to (A.12). And by Jensen's inequality (A.13) follows.

Remark 5. Though it is not necessary for our purposes, we point out

that Lemma 6 can be generalized considerably as follows:

If $Y_1,\ldots,Y_k$ are independent random variables with continuous distribution

functions $F_1,\ldots,F_k$, then under mild conditions on $F_1,\ldots,F_k$ we have

$$(A.15) \quad E(\max_{j=1,\ldots,k} Y_j \mid \max_{j=1,\ldots,k-1} Y_j = \xi) = \xi + \int_\xi^\infty [1 - F_k(\eta)] d\eta, \quad \xi \in \mathbb{R},$$

where the r.h.s. is convex in $\xi \in \mathbb{R}$. Thus Jensen's inequality implies

$$(A.16) \quad E(\max_{j=1,\ldots,k} Y_j) \geq E(\max_{j=1,\ldots,k-1} Y_j) + \int_{E(\max_{j=1,\ldots,k-1} Y_j)}^\infty [1 - F_k(\eta)] d\eta.$$

On the other extreme, for standard normal variables, we have

$$E(\max_{j=1,\ldots,k} V_j \mid \max_{j=1,\ldots,k-1} V_j) = T(\max_{j=1,\ldots,k-1} V_j) \quad \text{and}$$

$$E_k(\underline{0},\rho) = E(T(\max_{j=1,\ldots,k-1} V_j)) \geq T(E_{k-1}(\underline{0},\rho)) \ldots \geq \underbrace{T \circ \ldots \circ}_{k-1} T(0).$$

Lemma 7. $E_k(\underline{\mu},\rho) \geq \mu_{[k]} + A_\beta, \quad \beta = 3,4, \quad \underline{\text{where}}$

$$(A.17) \quad A_3 = \rho \, T(\rho^{-1}[a_{k-1} - k(k-1)^{-1}(\mu_{[k]} - \bar{\mu})])$$

$$\underline{\text{with}} \quad \bar{\mu} = k^{-1}(\mu_1 + \ldots + \mu_k), \quad \underline{\text{and}}$$

$$(A.18) \quad A_4 = \rho \, T\{\rho^{-1}(\mu_{[k-1]} - \mu_{[k]}) + 2^{1/2} T(2^{-1/2}\rho^{-1}(\mu_{[k-2]} - \mu_{[k-1]}))\}.$$

Proof: By (A.5)

$$E_k(\underline{\mu},\rho) = \mu_{[k]} + \int_{\mathbb{R}} \Phi(\rho^{-1}\xi)[1 - \prod_{i=1}^{k-1} \Phi(\rho^{-1}(\xi+\mu_{[k]} - \mu_{[i]}))]d\xi,$$

which by the log-concavity of $\Phi$ is bounded from below by

$$\mu_{[k]} + \int_{\mathbb{R}} \Phi(\rho^{-1}\xi)[1-\Phi(\rho^{-1}(\xi+\mu_{[k]} - (k-1)^{-1}(\mu_{[1]}+\ldots+\mu_{[k-1]})))^{k-1}]d\xi.$$

Integrating by parts this in turn equals to

$$\mu_{[k]} + \rho\int_{\mathbb{R}} T(\rho^{-1}\xi)d\{\Phi(\rho^{-1}(\xi+k(k-1)^{-1}(\mu_{[k]} - \bar{\mu})))^{k-1}\}$$

$$= \mu_{[k]} + \rho E(T(\rho^{-1}\{\max_{j=1,\ldots,k-1} V_j - k(k-1)^{-1}(\mu_{[k]}-\bar{\mu})\})),$$

and by Jensen's inequality this is bounded from below by (A.17). (A.18) follows immediately by applying (A.11) (for k-1 instead of k) to (A.13).

Remark 6. One can get more lower bounds of type (A.18) by iterating (A.13) N times (N = 1,2,...,k-2) before applying (A.11) (for k-N instead of k).

Note that among $\underset{\sim}{A}_2$, $\underset{\sim}{A}_3$ and $\underset{\sim}{A}_4$ (or among $\tilde{A}_1$ and $\tilde{A}_2$) none of the bounds is uniformly better than the others. One reason is that for every fixed $\eta \in \mathbb{R}$, $u\, T(u^{-1}\eta)$ is strictly increasing in u.

REFERENCES

[1] Bickel, P. J. and Yahav, J. A. (1977). On selecting a set of good populations. Statistical Decision Theory and Related Topics II. (Gupta, S. S. and Moore, D. S. eds.). Academic Press, 37-55.

[2] Blumenthal, S. and Cohen, A. (1968). Estimation of the larger of two normal means. J. Amer. Statist. Assoc. 63, 861-876.

[3] Bratcher, T. L. and Bhalla, P. (1974). On the properties of an optimal selection procedure. Comm. Statist. 3, 191-196.

[4] Chernoff, H. and Yahav, J. A. (1977). A subset selection problem employing a new criterion. Statistical Decision Theory and Related Topics II. (Gupta, S. S. and Moore, D. S. eds.). Academic Press, 93-119.

[5] David, H. A. (1970). Order Statistics. John Wiley, New York.

[6] Deely, J. J. and Gupta, S. S. (1968). On the properties of subset selection procedures. Sankhyā Ser. A 30, 37-50.

[7] Goel, P. K. and Rubin, H. (1977). On selecting a subset containing the best population - A Bayesian approach. Ann. Statist. 5, 969-983.

[8] Gupta, S. S. (1956). On a decision rule for a problem in ranking means. Mimeo. Series No. 150, Inst. of Statist., Univ. of North Carolina, Chapel Hill, North Carolina.

[9] Gupta, S. S. (1965). On some multiple decision (selection and ranking) rules. Technometrics 7, 225-245.

[10] Gupta, S. S. and Hsu, J. C. (1978). On the performance of some subset selection procedures. Communications in Statistics, Vol. B7, No. 6.

[11] Gupta, S. S. and Hsu, J. C. (1977). On the monotonicity of Bayes subset selection procedures. Book 4, Proceedings of the 41st Session of the International Statistical Institute held in New Delhi, December 1977, 208-211.

[12] Gupta, S. S. and Miescke, K. J. (1978). On subset selection procedures for ranking means of three normal populations. Mimeo. Series No. 78-19, Dept. of Statist., Purdue University, W. Lafayette, Indiana.

[13] Lehmann, E. L. (1957). A theory of some multiple decision problems I. Ann. Math. Statist. 28, 1-25.

[14] Seal, K. C. (1957). An optimum decision rule for ranking means of normal populations. Calcutta Statist. Assoc. Bull. 7, 131-150.

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER<br>Mimeograph Series #78-23 | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle)<br>Bayesian Subset Selection for Additive and Linear Loss Functions | | 5. TYPE OF REPORT & PERIOD COVERED<br>Technical |
| | | 6. PERFORMING ORG. REPORT NUMBER<br>Mimeo. Series #78-23 |
| 7. AUTHOR(s)<br>Klaus-J. Miescke | | 8. CONTRACT OR GRANT NUMBER(s)<br>ONR N00014-75-C-0455 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br>Purdue University<br>Department of Statistics<br>West Lafayette, IN 47907 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS |
| 11. CONTROLLING OFFICE NAME AND ADDRESS<br>Office of Naval Research<br>Washington, DC | | 12. REPORT DATE<br>October 1978 |
| | | 13. NUMBER OF PAGES<br>23 |
| 14. MONITORING AGENCY NAME & ADDRESS(if different from Controlling Office) | | 15. SECURITY CLASS. (of this report)<br>Unclassified |
| | | 15a. 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)

Subset Selection, Bayes Rules, Additive Linear Loss Function, Gupta's Maximum Means Procedure, Approximate Bayes Solutions, Expected Values of the Maximum.

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

Given $k$ independent samples $\underline{X}_1,\ldots,\underline{X}_k$ of common size $n$ from $k$ populations $\pi_1,\ldots,\pi_k$ with distributions $P_{\theta_i}$, $\theta_i \in \Omega \subseteq \mathbb{R}$, $i = 1,\ldots,k$, the problem is to select a non-empty subset $S(\underline{X}_1,\ldots,\underline{X}_k)$ from $\{\pi_1,\ldots,\pi_k\}$, which is associated with "good" (large) $\theta$-values. More precisely, we consider this problem from a Bayesian approach. (over)

By choosing additive $(L(\underline{\theta},s) = \sum_{i \in s} \ell_i(\underline{\theta}))$ and especially linear $(L(\underline{\theta},s) = \sum_{i \in s} (\theta_{[k]} - \theta_i - \varepsilon))$ loss functions we try to fill a gap lying in between the results of Deely and Gupta (1968) and more recent papers due to Goel and Rubin (1977), Gupta and Hsu (1978) and other authors. It is shown that under a certain "normal model" Seal's procedure turns out to be Bayes w.r.t. an unrealistic loss function whereas Gupta's maximum means procedure turns out to be (for large n) asymptotically Bayes w.r.t. more realistic additive loss functions. Finally, in the appendix some bounds for $E(\max_{j=1,\ldots,k} (\mu_j + \rho V_j))$ are derived (where $\underline{\mu} \in \mathbb{R}^k$, $\rho \in \mathbb{R}$ are fixed known and $\underline{V} \sim N(\underline{0},I)$) to approximate the Bayes rules w.r.t. linear loss functions in cases where n is finite.