# BAYESIAN LOOK AHEAD ONE STAGE SAMPLING ALLOCATIONS FOR SELECTING THE LARGEST OF $k \geq 3$ NORMAL MEANS*

by

Shanti S. Gupta      and    Klaus J. Miescke
Purdue University            University of Illinois at Chicago

Technical Report #93-32

Department of Statistics
Purdue University

July 1993

---

o

# BAYESIAN LOOK AHEAD ONE STAGE SAMPLING ALLOCATIONS
## FOR SELECTING THE LARGEST OF $k \geq 3$ NORMAL MEAN*

by

| Shanti S. Gupta | and | Klaus J. Miescke |
|---|---|---|
| Purdue University | | University of Illinois at Chicago |
| West Lafayette, IN | | Chicago, IL |

## ABSTRACT

From $k \geq 3$ independent normal populations with unknown means and a common known variance, samples of unequal sizes are observed at stage 1. The goal is to find that population with the largest mean. Using the Bayes approach, optimum allocations of $m$ additional observations, at stage 2, are derived under the linear loss function.

AMS Subject Classification: Primary 62F07; secondary 62F15.

Key Words: Bayesian look ahead procedure; selection and ranking; normal populations.

# 1. Introduction

Let $\pi_1, \ldots, \pi_k$ be $k \geq 3$ given normal populations with unknown means $\theta_1, \ldots, \theta_k \in \mathbf{R}$ and a common known variance $\sigma^2 > 0$. Suppose we want to find that population with the largest mean using a Bayes selection rule which is based on a known prior density $\pi(\boldsymbol{\theta})$ for $\boldsymbol{\theta}$, and a given loss function $L(\boldsymbol{\theta}, i)$ for selecting $i \in \{1, 2, \ldots, k\}$ at $\boldsymbol{\theta} \in \mathbf{R}^k$. The loss function adopted later on to derive specific results will be linear, i.e. $L(\boldsymbol{\theta}, i) = \theta_{[k]} - \theta_i$, $i = 1, \ldots, k$, $\boldsymbol{\theta} \in \mathbf{R}^k$, where $\theta_{[k]} = \max\{\theta_1, \ldots, \theta_k\}$. Assume now that we are in the situation where $k$ independent samples of sizes $n_1, \ldots, n_k$, respectively, have been observed already at a first stage, and where $m$ additional observations are still allowed to be taken at a future second stage. The problem of interest treated below is how to allocate these $m$ observations in an optimum way among the $k$ populations, given all the information, prior and first stage observations, gathered so far. It should be pointed out that the special case of $n_1 = \ldots = n_k = 0$ represents the analogous problem of how to allocate $m$ observations at a first stage. Thus, the results derived below will provide also optimum sampling plans for stage 1.

Looking ahead one stage using the expected posterior Bayes risk, given all observations collected so far, does not only lead the way to an optimum allocation of observations in the future. It also provides in a snapshot a relative measure of how much better the decision can be expected to be after further sampling has been performed following this optimum allocation. If costs of sampling at the second stage would be incorporated in the loss, which is not done here, an optimum stopping rule could be implemented, too. In many empirical studies in marketing research (e.g. direct marketing), medical research (e.g. clinical trials, cf. Whitehead (1991)), and social research (e.g. survey sampling), there are interim analyses performed at certain stages to decide if sampling should be continued, and if so, how to allocate observations.

Under the assumption of $k$ independent normal priors and either a linear loss or a 0-1 loss, a solution to the problem has been obtained previously by Gupta and Miescke (1993) for the case of $k = 2$, which turns out to be already rather involved. It allocates observations in such a way that the posterior gets as close as possible to being decreasing in transposition (DT). Moreover, somewhat surprising, it does not depend at all on the

2

observations gathered at the first stage. This fact implies that one can allocate, in an optimum way, one new observation at a time, until all $m$ have been drawn, thereby arriving at the same allocation as in the former approach. For $k \geq 3$ populations, this is no longer true! As it will be shown below, the values observed at the first stage are in fact relevant for the allocation of further observations.

In a fully sequential approach, the optimum allocation could be constructed, at least in principle, by finding first the Bayesian terminal selection rules based on all possible $m$ allocations and then use backward induction to optimize successively the allocations before. Although the former are not hard to find, the latter appears to be infeasible to carry out in practice. Therefore, the procedure proposed here is to optimally allocate one observation at a time, pretending that it is the last to be drawn, and then iterate this process until $m$ observations have been taken. Allocating more than one observation at a time, on the other hand, appears to be less appealing since with each new observation more is learnt about the unknown parameters, which improves the basis for further decisions. Look ahead procedures in various settings are described and discussed in Berger (1985).

Selecting the population with the largest (overall) sample mean is usually called the natural selection rule, since it is the uniformly best permutation invariant selection procedure, in the frequentist sense, for a general class of loss functions, provided that the sample sizes are equal. However, for unequal sample sizes, the natural selection rule loses much of its quality, although it still remains intuitively appealing. Therefore, optimum sample size allocations for this rule have been considered in the frequentist approach by Bechhofer (1969), Dudewicz and Dalal (1975), and Bechhofer, Hayter, and Tamhane (1991). On the other hand, Bayes rules with normal priors turn out to have complicated forms which cannot be represented in closed form, except for those situations where the posterior is (DT), as it has been shown in Gupta and Miescke (1988). Bayes rules for similar but more involved models have been studied by Berger and Deely (1988) and Fong and Berger (1993). For the present setting, it was recommended in Gupta and Miescke (1988) to plan the experiment's sampling allocation in such a way as to make the posterior (DT). This is corroborated by the optimum allocation for $k = 2$. However, for more than two populations, it will be shown below that there is some tradeoff needed between this recommendation and the values observed at the first stage.

## 2. General Framework and Notation

For convenience and to facilitate comparisons with Gupta and Miescke (1993), the same notation used there will also be adopted here. After a standard reduction of the data by sufficiency, the model assumptions can be summarized as follows: At $\theta = (\theta_1, \ldots, \theta_k) \in \mathbf{R}^k$, $X_i \sim N(\theta_i, p_i^{-1})$ with $p_i^{-1} = \sigma^2/n_i$, and $Y_i \sim N(\theta_i, q_i^{-1})$ with $q_i^{-1} = \sigma^2/m_i$, are the sample means from the samples of population $\pi_i$ at stage 1 and stage 2, respectively, $i = 1, \ldots, k$, which altogether are assumed to be independent. A priori, the parameters $\underline{\Theta} = (\Theta_1, \ldots, \Theta_k)$ are considered to be random and follow a given prior distribution which will be later assumed to satisfy $\Theta_i \sim N(\mu_i, \nu_i^{-1})$, $i = 1, \ldots, k$, and are independent. Let the loss function for selections at stage 1 and stage 2 be denoted by $L(\theta, s)$, $\theta \in \mathbf{R}^k$, $s = 1, \ldots, k$, which will be later assumed to be linear. Costs of sampling are not incorporated in the loss to simplify the presentation of basic ideas. The stopping rule is considered here to be deterministic: the sampling process stops after $m$ observations are taken.

Having observed $\underline{X} = \underline{x} \in \mathbf{R}^k$ at stage 1, every Bayes selection rule $d_1^*(\underline{x})$ satisfies

$$E\{L(\underline{\Theta}, d_1^*(\underline{x}))|\underline{X} = \underline{x}\} = \min_{i=1,\ldots,k} E\{L(\underline{\Theta}, i)|\underline{X} = \underline{x}\}. \tag{1}$$

Likewise, after $\underline{Y} = \underline{y} \in \mathbf{R}^k$ has been observed at stage 2, every Bayes rule satisfies

$$E\{L(\underline{\Theta}, d_2^*(\underline{x}, \underline{y}))|\underline{X} = \underline{x}, \underline{Y} = \underline{y}\} = \min_{i=1,\ldots,k} E\{L(\underline{\Theta}, i)|\underline{X} = \underline{x}, \underline{Y} = \underline{y}\}. \tag{2}$$

There will be no need to consider randomized Bayes selection rules in the following since minimaxity and invariance concepts will not be used.

Many results for Bayes selection rules can be found in the literature. An overview is provided by Gupta and Panchapakesan (1979, 1991). Only recently, however, attention has been given also to non-symmetric models. The binomial case has been treated by Abughalous and Miescke (1989), whereas the normal case is studied and discussed in Gupta and Miescke (1988), and more involved models are treated in Berger and Deely (1988) and in Fong and Berger (1993). Rather than studying the properties of the Bayes selection rules $d_1^*$ and $d_2^*$ in details, let us assume here that they have been derived already, are ready to be used, and that all that is left to do is to allocate sample sizes in an optimum manner.

4

Let us first consider fixed total sample size allocation problems. Before entering each of the two stages, similar allocation problems of this type arise, which are closely related. Before entering stage 1, by looking ahead one stage, we would like to minimize the expected posterior risk subject to $n_1 + \ldots + n_k = n$, where $n$ is the total number of observations allowed to be taken at stage 1. This leads to the following criterion for $n_1, \ldots, n_k$

$$\min_{\substack{n_1, \ldots, n_k \\ n_1 + \ldots + n_k = n}} E(\min_{i=1, \ldots, k} E\{L(\underline{\Theta}, i) | \underline{X}\}). \tag{3}$$

Likewise, at the end of stage 1, the criterion for $m_1, \ldots, m_k$ with a total number of $m$ observations allowed at stage 2 is the following

$$\min_{\substack{m_1, \ldots, m_k \\ m_1 + \ldots + m_k = m}} E\{\min_{i=1, \ldots, k} E\{L(\underline{\Theta}, i) | \underline{X} = \underline{x}, \underline{Y}\} | \underline{X} = \underline{x}\}. \tag{4}$$

As explained in Gupta and Miescke (1993), one can get solutions of (3) from those of (4) by setting $n_1 = \ldots = n_k = 0$ and then relabel $m_i$ by $n_i$, $i = 1, \ldots, k$, and $m$ by $n$. Thus, we need to consider only criterion (4) in the sequel, since it is more general.

In the following, let us assume independent normal priors, i.e. $\Theta_i \sim N(\mu_i, \nu_i^{-1})$, $i = 1, \ldots, k$, which are independent. To evaluate the inner conditional expectation of (4), note that the conditional distribution of $\underline{\Theta}$, given $\underline{X} = \underline{x}$ and $\underline{Y} = \underline{y}$ is:

$$\text{Given } \underline{X} = \underline{x} \text{ and } \underline{Y} = \underline{y}, \quad \Theta_1, \ldots, \Theta_k \text{ are independent with} \tag{5}$$
$$\Theta_i \sim N\left(\frac{\alpha_i \mu_i(\underline{x}) + q_i y_i}{\alpha_i + q_i}, \frac{1}{\alpha_i + q_i}\right), \quad i = 1, \ldots, k,$$

where $\alpha_i = p_i + \nu_i$ and $\mu_i(\underline{x}) = (\nu_i \mu_i + p_i x_i)/(\nu_i + p_i)$, $i = 1, \ldots, k$.

To evaluate the outer conditional expectation in (4), the conditional distribution of $\underline{Y}$, given $\underline{X} = \underline{x}$, has to be used, which is:

$$\text{Given } \underline{X} = \underline{x}, \quad Y_1, \ldots, Y_k \text{ are independent with} \tag{6}$$
$$Y_i \sim N\left(\mu_i(\underline{x}), \frac{\alpha_i + q_i}{\alpha_i q_i}\right), \quad i = 1, \ldots, k.$$

The posterior distribution of $\Theta_1, \ldots, \Theta_k$ at stage 2, as given by (5), is (DT) if and only if

$$\alpha_1 + q_1 = \alpha_2 + q_2 = \ldots = \alpha_k + q_k. \tag{7}$$

5

It was recommended by Gupta and Miescke (1988) to plan the experiment in such a way that (7) is satisfied, because of two reasons. First, the Bayes rule is then of a very simple closed form, and second, the posterior information about the $k$ unknown parameters is then equally and fairly balanced. The solutions of criterion (4) for $k = 2$ under linear and 0-1 loss, as shown in Gupta and Miescke (1993), turn out to be the same: Choose $m_1$ and $m_2$ subject to $m_1 + m_2 = m$ in such a way that one gets as close as possible to the (DT) configuration (7). Since this common solution does not depend at all on the observations $\underline{X} = \underline{x}$ at stage 1, the solution has been found also to the analogous problem where one observation at a time, pretending it is the last to be observed, has to be allocated successively. Using the latter $m$ times leads to the same result as performing the former in one single step. Both procedures are in the case of $k = 2$ identical with the optimum truncated sequential allocation rule which employs Bayes terminal selections and optimum allocations determined by backward induction. As will be seen later, this is no longer true for the case of $k \geq 3$.

Taking into account that for $k \geq 3$, optimum future sampling allocations do depend on the past observations, it seems more appropriate to break down the allocation of $m$ observations into more than one step. Several sampling allocation schemes are reasonable. To represent them in a convenient way, let us introduce some useful notation. For $t \leq m$, let $\mathcal{R}_t$ denote the fixed sample size $t$ allocation determined by (4) with $m$ replaced by $t$ there. Furthermore, let $\mathcal{R}_{t,1}$ allocate any one, but only one, observation to one of the $l \leq t$ populations suggested by $\mathcal{R}_t$. Finally, let $\mathcal{B}_1$ stand for the optimum allocation of one observation, determined by known future actions. The procedure which allocates the last $m-1$ observations using criterion (4), with $m$ replaced by $m-1$ there, and allocates the first observation through backward optimization, can thus be represented by $(\mathcal{B}_1, \mathcal{R}_{m-1})$. In a similar manner, $(\mathcal{B}_1, \mathcal{R}_{m-1,1}, \mathcal{R}_{m-2})$ allocates the last $m-2$ observation using criterion (4), with $m$ replaces by $m-2$ there, allocates the second observation to one of the populations suggested by $\mathcal{R}_{m-1}$ after the first observation has been drawn, and finally allocates the first observation by backward optimization. It should be pointed out that $\mathcal{R}_t$ and $\mathcal{R}_{t,1}$ are stand-alone procedures, whereas $\mathcal{B}_1$ is only meaningful in connection with completely specified actions in the future.

The best sampling allocation of $m$ observations is determined by an optimum allocation of the $m - th$ observation, and backward induction. This is $(\mathcal{B}_1, \mathcal{B}_1, \ldots, \mathcal{B}_1, \mathcal{R}_1)$, with $m - 1$ repetitions of $\mathcal{B}_1$. For comparisons of other allocation schemes, let $\mathcal{S} \prec (=) \mathcal{T}$ mean that procedure $\mathcal{T}$ has a smaller (the same) overall Bayes risk than procedure $\mathcal{S}$. Using the basic fact that in the present situation every relevant (conditional) expectation of a maximum is larger than the maximum of the respective (conditional) expectations, the chains of preferences presented below in the two lemmas can be verified easily. Therefore, their proofs are omitted for brevity.

**Lemma 1.**

$$(\mathcal{R}_m) \prec (\mathcal{R}_{m,1}, \mathcal{R}_{m-1}) \prec (\mathcal{B}_1, \mathcal{R}_{m-1}) \prec (\mathcal{B}_1, \mathcal{R}_{m-1,1}, \mathcal{R}_{m-2}) \prec \qquad (8)$$
$$\prec (\mathcal{B}_1, \mathcal{B}_1, \mathcal{R}_{m-2}) \prec \ldots \prec (\mathcal{B}_1, \mathcal{B}_1, \ldots, \mathcal{B}_1, \mathcal{R}_1).$$

Since backward optimization, i.e. using $\mathcal{B}_1$, appears to be infeasible, one may be inclined to capitalize on one of the following sampling allocation rules:

**Lemma 2.**

$$(\mathcal{R}_m) \prec (\mathcal{R}_{m,1}, \mathcal{R}_{m-1}) \prec (\mathcal{R}_{m,1}, \mathcal{R}_{m-1,1}, \mathcal{R}_{m-2}) \prec \ldots \prec \qquad (9)$$
$$\prec (\mathcal{R}_{m,1}, \mathcal{R}_{m-1,1}, \ldots, \mathcal{R}_{2,1}, \mathcal{R}_1).$$

There is one problem that arises when one attempts to use procedure $\mathcal{R}_{t,1}$ : it is not clear from which one of the $l \leq t$ populations, suggested by $\mathcal{R}_t$, should the next observation be drawn in an optimum fashion. The natural alternative is to consider using one of the following simpler allocation rules:

$$(\mathcal{R}_1, \mathcal{R}_{m-1}), \ (\mathcal{R}_1, \mathcal{R}_1, \mathcal{R}_{m-2}), \ \ldots, \ (\mathcal{R}_1, \mathcal{R}_1, \ldots, \mathcal{R}_1). \qquad (10)$$

Also here some conflict arises which, however, appears to be less serious. Whenever an $\mathcal{R}_{t,1}$ is replaced by an $\mathcal{R}_1$, it remains unclear whether the latter picks one population from those available for $\mathcal{R}_{t,1}$ from $\mathcal{R}_t$, i.e. whether $\mathcal{R}_1$ is a version of $\mathcal{R}_{t,1}$.

In the next section, sampling rules $\mathcal{R}_1$ and $\mathcal{R}_m$ will be studied in more detail under the linear loss function. The results for $\mathcal{R}_m$ can also be applied to $\mathcal{R}_t$ for $t \leq m$ in an obvious manner.

7

## 3. Main Results For Linear Loss

In this section we assume that the loss in linear, i.e. $L(\boldsymbol{\theta}, i) = \theta_{[k]} - \theta_i$, $i = 1, \ldots, k$, $\boldsymbol{\theta} \in \mathbf{R}^k$, where $\theta_{[k]} = \max\{\theta_1, \ldots, \theta_k\}$. Criterion (4) reduces then to minimize as a function of $m_1, \ldots, m_k$, subject to $m_1 + \ldots + m_k = m$, the look ahead expected posterior risk

$$E\{\Theta_{[k]} | \underline{X} = \underline{x}\} - E\{\max_{i=1,\ldots,k} E\{\Theta_i | \underline{X} = \underline{x}, \underline{Y}\} | \underline{X} = \underline{x}\}. \tag{11}$$

For $k \geq 2$ and the independent normal priors introduced in section 2, this reduces further by using first (5) and then (6) to maximize as a function of $m_1, \ldots, m_k$, subject to $m_1 + \ldots + m_k = m$, the look ahead expected gain

$$E\{\max_{i=1,\ldots,k} \frac{\alpha_i \mu_i(\underline{x}) + q_i Y_i}{\alpha_i + q_i} | \underline{X} = \underline{x}\} = E\left(\max_{i=1,\ldots,k} \left[\mu_i(\underline{x}) + \left(\frac{q_i}{\alpha_i(\alpha_i + q_i)}\right)^{\frac{1}{2}} N_i\right]\right), \tag{12}$$

where $N_1, \ldots, N_k$ are independent standard normal generic random variables.

To study the properties $E(max_{i=1,\ldots,k}[a_i + b_i N_i])$ as a function of $a_i \in \mathbf{R}$ and $b_i > 0$, $i = 1, \ldots, k$, it proves useful to introduce the following auxiliary function T, given by

$$T(w) = w \, \Phi(w) + \varphi(w) = \int_{-\infty}^{w} \Phi(u) \, du \, , \, w \in \mathbf{R}, \tag{13}$$

where $\Phi$ and $\varphi$ denote the standard normal c.d.f. and density, respectively. This function has been introduced and studied previously in Miescke (1979), but only for the case where $b_1, \ldots, b_k$ are equal. T is positive, strictly increasing, and convex. Moreover, it has the following basic properties, which can be verified easily: $T(w) > w$, for $w > 0$; $T(w) = T(-w) + w$, for $w \in \mathbf{R}$; $\gamma \, T(w/\gamma)$ is increasing in $\gamma \in \mathbf{R}$ for $w \in \mathbf{R}$; $T(w) = E(max[N_1, w])$, for $w \in \mathbf{R}$. To simplify notation, let $M_k = max_{i=1,\ldots,k}[a_i + b_i N_i]$ with $a_i \in \mathbf{R}$ and $b_i > 0$, $i = 1, \ldots, k$, in the following.

## Lemma 3.

(i) $E(T(M_1)) = (b_1^2 + 1)^{1/2} \, T(a_1/(b_1^2 + 1)^{1/2})$.

(ii) $E(M_2) = a_1 + (b_1^2 + b_2^2)^{1/2} \, T((a_2 - a_1)/(b_1^2 + b_2^2)^{1/2})$.

(iii) $E(M_k) = a_k + b_k E(T((M_{k-1} - a_k)/b_k) > a_k + b_k T((E(M_{k-1}) - a_k)/b_k)$.

(iv) $E(M_k)$ *is increasing in* $a_i \in \mathbf{R}$ *and* $b_i > 0$, $i = 1, \ldots, k$.

**Proof:** (i) Let $a_1 = a$ and $b_1 = b$ for brevity. Then one can see that

$$E(T(M_1)) = a \int_{\mathbf{R}} \Phi(a + bz) \; \varphi(z) \; dz + b \int_{\mathbf{R}} \Phi(a + bz) \; z \; \varphi(z) \; dz$$

$$+ \int_{\mathbf{R}} \varphi(a + bz) \; \varphi(z) \; dz$$

$$= a \; \Phi(a/(b^2 + 1)^{1/2}) + (b^2 + 1) \int_{\mathbf{R}} \varphi(a + bz) \; \varphi(z) \; dz.$$

Using integration by parts, the second integral could be brought into the form of the third. Combining now the two $\varphi$-functions into one exponential function, standard calculations show that (i) holds.

(iii) From the representation $M_k = a_k + b_k E(max[N_k, (M_{k-1} - a_k)/b_k])$ and the identity $T(w) \equiv E(max[N_1, w])$, one can see that the conditional expectation of $M_k$ , given $N_1, \ldots, N_{k-1}$ , is equal to $a_k + b_k T((M_{k-1} - a_k)/b_k)$. The inequality follows from Jensen's inequality applied to the convex function $T$.

(ii) This is seen by applying (i) to the equation in (iii) with $k = 2$.

(iv) The monotonicity with respect to $a_1, \ldots, a_k$ is obviously correct. As to $b_1, \ldots, b_k$ , consider without loss of generality $b_k$. We know that $\gamma \; T(w/\gamma)$ is increasing in $\gamma$ for every $w \in \mathbf{R}$. An application of this fact to the equation in (iii) completes the proof.

**Remark 1.** The equation and the inequality in (iii) of the lemma can be used iteratively to get lower bounds on $E(M_k)$ in terms of $E(M_r)$ for $r = 2, \ldots, k - 1$. $E(M_2)$ has a closed form representation given by (ii) in the lemma.

The remainder of this section is devoted to a detailed study of the properties of the allocation rule $\mathcal{R}_1$, and finally to some discussion of optimum allocations for stage 1. $\mathcal{R}_1$ is found by maximizing (12), subject to one of the values $q_1, \ldots, q_k$ being equal to q and all others being equal to 0, where $q = 1/\sigma^2$. All of the results derived below are given in terms of $q$, and it should be pointed out that they also hold true for any $q = h/\sigma^2$ with $1 \leq h \leq m$, i.e. for the analogous allocation rule which optimizes allocation of $h$ observations to exactly one population, where $h \in \{1, \ldots, m\}$. For this case, criterion (12) simplifies to a closed form as given in the following lemma, which will serve as the main tool to set up a useful algorithm for finding the desired optimum sampling allocation.

**Lemma 4.** *For $q_i = q$ , and for $q_j = 0$ , $j \neq i$ , we have*

$$E\{ \max_{s=1,\ldots,k} \frac{\alpha_s \mu_s(\underline{x}) + q_s Y_s}{\alpha_s + q_s} | \underline{X} = \underline{x} \} = \mu_i(\underline{x}) + \sigma_i T(\Delta_i / \sigma_i) , \qquad (14)$$

*where $\sigma_i = \left( \frac{q}{\alpha_i(\alpha_i + q)} \right)^{\frac{1}{2}}$ , and $\Delta_i = max_{j \neq i}[\mu_j(\underline{x}) - \mu_i(\underline{x})]$.*

**Proof:** Under the assumptions above, the right hand side of (12) can be written as $\mu_i(\underline{x}) + \sigma_i E(max[N_i, \Delta_i / \sigma_i])$. The rest follows from the identity $T(w) \equiv E(max[N_i, w])$.

In the sequel, let $\mu_{[1]}(\underline{x}) < \mu_{[2]}(\underline{x}) < \ldots < \mu_{[k]}(\underline{x})$ denote the ordered posterior means from stage 1. Let $\pi_{(i)}$ be the population which is associated with $\mu_{[i]}(\underline{x})$, and let $\alpha_{(i)}$, $\sigma_{(i)}$, and $\Delta_{(i)}$ be defined analogously. Furthermore, let $\mathcal{R}^{(i)}(\underline{x})$ , in short $\mathcal{R}^{(i)}$ , denote the rule which assigns, at $\underline{X} = \underline{x}$ , the next allocation to population $\pi_{(i)}$ , $i = 1, \ldots, k$.

First, let us consider the two populations $\pi_{(k-1)}$ and $\pi_{(k)}$ , since they turn out to play a very special role in this situation: these are the only two population between which a preference in terms of the order relation " $\prec$ " can be established which does not depend on $\mu_1(\underline{x}), \ldots, \mu_k(\underline{x})$. This fact is established in the following theorem.

**Theorem 1.**

$$\alpha_{(k-1)} > (=, <) \alpha_{(k)} \quad \text{if and only if} \quad \mathcal{R}^{(k-1)} \prec (=, \succ) \mathcal{R}^{(k)} .$$

**Proof:** Consider the look ahead expected gain, given $\underline{X} = \underline{x}$ , of rule $\mathcal{R}^{(k)}$ , which is based on one more observation from population $\pi_{(k)}$. Let it be denoted by $g_k(\underline{x})$. Starting with its representation given by (14), we can see that it satisfies

$$g_k(\underline{x}) = \mu_{[k]}(\underline{x}) + \sigma_{(k)} T(\Delta_{(k)} / \sigma_{(k)}) = \mu_{[k-1]}(\underline{x}) + \sigma_{(k)} T(\Delta_{(k-1)} / \sigma_{(k)}) , \qquad (15)$$

where the second equation follows from $\Delta_{(k)} = \mu_{[k-1]} - \mu_{[k]} = -\Delta_{(k-1)}$ , and from the identity $T(w) \equiv T(-w) + w$. On the other hand, the look ahead expected gain $g_{k-1}(\underline{x})$ of rule $\mathcal{R}^{(k-1)}$ is seen to be

$$g_{k-1}(\underline{x}) = \mu_{[k-1]} + \sigma_{(k-1)} T(\Delta_{(k-1)} / \sigma_{(k-1)}). \qquad (16)$$

The rest follows from the fact that $\sigma T(\Delta / \sigma)$ is increasing in $\sigma$ for every $\Delta \in \mathbf{R}$.

Let $\mathcal{R}^{(*)}(\underline{x})$ be the better of the two rules $\mathcal{R}^{(k-1)}$ and $\mathcal{R}^{(k)}$ at $\underline{X} = \underline{x}$, as determined by theorem 1, and let $\alpha_{(*)} = min[\alpha_{(k-1)}, \alpha_{(k)}]$ and $\sigma_{(*)} = max[\sigma_{(k-1)}, \sigma_{(k)}]$ be the parameters associated with that population $\pi_{(*)}$, say, which is chosen by $\mathcal{R}^{(*)}(\underline{x})$. This rule has the look ahead expected gain

$$g_*(\underline{x}) = \mu_{[k-1]} + \sigma_{(*)}T((\mu_{[k]} - \mu_{[k-1]})/\sigma_{(*)}) , \qquad (17)$$

which has to be compared now with the maximum look ahead expected gain of the rules $\mathcal{R}^{(1)}, \ldots, \mathcal{R}^{(k-2)}$, i.e. with the maximum of

$$g_t(\underline{x}) = \mu_{[t]} + \sigma_{(t)}T((\mu_{[k]} - \mu_{[t]})/\sigma_{(t)}) , \ t = 1, \ldots, k-2. \qquad (18)$$

Comparing any two rules $\mathcal{R}^{(i)}$ and $\mathcal{R}^{(j)}$ with $1 \leq i < j \leq k-2$ will necessarily involve $\mu_{[i]}(\underline{x})$, $\mu_{[j]}(\underline{x})$, and $\mu_{[k]}(\underline{x})$. This situation is summarized in the next theorem.

**Theorem 2.** *For $1 \leq i < j \leq k-2$ the following holds.*

**(a)** *If $\alpha_{(i)} \geq \alpha_{(j)}$, i.e. if $\sigma_{(i)} \leq \sigma_{(j)}$, then $\mathcal{R}^{(i)} \prec \mathcal{R}^{(j)}$.*

**(b)** *If $\alpha_{(i)} < \alpha_{(j)}$, and $\sigma_{(i)} < (=,>) \sigma_{ij}$ then $\mathcal{R}^{(i)} \prec (=,\succ) \mathcal{R}^{(j)}$,*

*where $\sigma_{ij}$ is determined by*

$$\mu_{[j]} + \sigma_{(j)}T((\mu_{[k]} - \mu_{[j]})/\sigma_{(j)}) = \mu_{[i]} + \sigma_{ij}T((\mu_{[k]} - \mu_{[i]})/\sigma_{ij}) . \qquad (19)$$

*Moreover, the threshold $\sigma_{ij}$ for $\sigma_{(i)}$ satisfies*

$$\sigma_{(j)} < \sigma_{ij} < \frac{\mu_{[k]} - \mu_{[i]}}{\mu_{[k]} - \mu_{[j]}} \sigma_{(j)} . \qquad (20)$$

**Proof:** Let $i$ and $j$ satisfy the assumptions of the theorem. Note, that in particular we have $\mu_{[i]} < \mu_{[j]} < \mu_{[k-1]} < \mu_{[k]}$.

(a) Suppose that $\alpha_{(i)} \geq \alpha_{(j)}$. Then the following can be shown to hold true:

$$\begin{aligned}
g_i(\underline{x}) &= \mu_{[i]} + \sigma_{(i)}T(([\mu_{[k]} - \mu_{[j]}] + [\mu_{[j]} - \mu_{[i]}])/\sigma_{(i)}) \\
&< \mu_{[i]} + (\mu_{[j]} - \mu_{[i]}) + \sigma_{(i)}T((\mu_{[k]} - \mu_{[j]})/\sigma_{(i)}) \\
&= \mu_{[j]} + \sigma_{(i)}T((\mu_{[k]} - \mu_{[j]})/\sigma_{(i)}) \\
&\leq \mu_{[j]} + \sigma_{(j)}T((\mu_{[k]} - \mu_{[j]})/\sigma_{(j)}) = g_j(\underline{x}).
\end{aligned}$$

11

The first inequality follows from the fact that

$$T(u+v) - T(u) = \int_u^{u+v} \Phi(z) \, dz < v \, , \quad \text{for} \ \ v > 0 \ \ \text{and} \ \ u \in \mathbf{R}.$$

The second inequality follows from $\sigma_{(i)} \leq \sigma_{(j)}$ and from the fact that $\sigma T(\Delta/\sigma)$ is increasing in $\sigma$. This completes the proof of (a).

**(b)** Suppose now that $\alpha_{(i)} < \alpha_{(j)}$ , and consider the following difference as a function of $\sigma_{(i)}$ :

$$g_j(\underline{x}) - g_i(\underline{x}) = \left[\mu_{[j]} + \sigma_{(j)} T((\mu_{[k]} - \mu_{[j]})/\sigma_{(j)})\right] - \left[\mu_{[i]} + \sigma_{(i)} T((\mu_{[k]} - \mu_{[i]})/\sigma_{(i)})\right] \, , \quad (21)$$

which is positive (negative) if $\sigma_{(i)}$ tends to $\sigma_{(j)}$ (infinity) , and it is decreasing in $\sigma_{(i)}$ . Therefore, $\sigma_{ij}$ is determined as claimed in part (b) of the theorem.

Finally, we have to show that the threshold $\sigma_{ij}$ for $\sigma_{(i)}$ is located in the range given at the end of the theorem. As mentioned above, if $\sigma_{(i)}$ tends to $\sigma_{(j)}$ then (21) is positive, ant thus $\sigma_{ij} > \sigma_{(j)}$. On the other hand, if $\sigma_{(i)} = \sigma_{(j)}(\mu_{[k]} - \mu_{[i]})/(\mu_{[k]} - \mu_{[j]})$ then (21) reduces to

$$g_j(\underline{x}) - g_i(\underline{x}) = \mu_{[j]} - \mu_{[i]} - \frac{\mu_{[j]} - \mu_{[i]}}{\mu_{[k]} - \mu_{[j]}} \, \sigma_{(j)} \, T((\mu_{[k]} - \mu_{[j]})/\sigma_{(j)}) \, ,$$

which is negative, since $T(w) > w$ for $w \in \mathbf{R}$. This completes the proof of the theorem.

Let $\mathcal{R}^{(\bullet)}(\underline{x})$ be the best of the rules $\mathcal{R}^{(t)}$ for $t = 1, \ldots, k-2$ at $\underline{X} = \underline{x}$ , as determined by theorem 2, which has to be compared now with $\mathcal{R}^{(*)}(\underline{x})$. Let $\alpha_{(\bullet)}$, $\sigma_{(\bullet)}$, and $\mu_{[\bullet]}$ be the parameters associated with that population chosen by $\mathcal{R}^{(\bullet)}(\underline{x})$. This rule has the look ahead expected gain

$$g_\bullet(\underline{x}) = \mu_{[\bullet]} + \sigma_{(\bullet)} T((\mu_{[k]} - \mu_{[\bullet]})/\sigma_{(\bullet)}) \, . \quad (22)$$

This has to be compared now with (17). The result is summarized in the following

**Corollary.** *The comparison of the look ahead expected gains of procedures $\mathcal{R}^{(\bullet)}$ and $\mathcal{R}^{(*)}$ is the same as (a) and (b) in theorem 2, with $(i)$, $(j)$, $\sigma_{ij}$, $\mu_{[i]}$, and $\mu_{[j]}$ replaced by $(\bullet)$, $(*)$, $\sigma_{\bullet *}$, $\mu_{[\bullet]}$, and $\mu_{[k-1]}$, respectively.*

**Remark 2.** To apply allocation rule $\mathcal{R}_1$ at $\underline{X} = \underline{x}$, a quick and simple algorithm is to just find the maximum of (15), (16), and (18), and then to take the next observation from that population which has yielded the maximum. The results presented in the two theorems and in the corollary are not needed for this purpose.

**Remark 3.** The importance of the two theorems and the corollary is the insight which they provide into how and why the $k$ populations are ordered with respect to " $\prec$ ". Interpreting $\alpha_i = \nu_i + n_i/\sigma^2$ as the total of prior plus first stage sampling information gathered from population $\pi_i$ , $i = 1, \ldots, k$, we can see that the allocation rule $\mathcal{R}_1$ favors populations which have, at stage 1, larger yields $\mu_{[i]}$ and smaller sampling informations $\alpha_i$ , $i = 1, \ldots, k$. One exception, treated by theorem 1, is that for the two populations with the largest yields $\mu_{[k-1]}$ and $\mu_{[k]}$, that one with the smaller of the two values $\alpha_{(k-1)}$ and $\alpha_{(k)}$ is preferred. Besides this exception, however, both goals cannot always be met simultaneously: there is a tradeoff between the two goals which is provided by theorem 2 and by the corollary.

To conclude this study, it is natural to consider also the optimization of sampling allocations for stage 1. As mentioned below of (4), the results derived in this section are more general and apply also to this case. All that has to be done is to set $n_1 = n_2 = \ldots = n_k = 0$ and then relabel $m_i$ by $n_i$, for $i = 1, \ldots, k$, and $m$ by $n$. The relevant parameters reduce then to $\alpha_i = \nu_i$ and $\mu_i(\underline{x}) = \mu_i$, $i = 1, \ldots, k$, and the results for allocation rule $\mathcal{R}_1$ derived above hold analogously.

One special case, however, deserves further consideration: the case of *i.i.d.* priors, i.e. where $\Theta_i \sim N(\mu, \nu^{-1})$, $i = 1, \ldots, k$, and are independent. In this situation, (14) reduces to $\mu + \sigma_i/\sqrt{2\pi}$, and thus rule $\mathcal{R}_1$ allocates one observation to that population $\pi_i$ which has the smallest $\alpha_i$, $i = 1, \ldots, k$.

Suppose that an allocation of $m$ observations can be made in such a way that $\nu_1 + n_1/\sigma^2 = \ldots = \nu_k + n_k/\sigma^2$ holds, i.e. that the posterior at stage 1 is (DT). Then allocating $n$ times iteratively one observation to that population with the smallest $\alpha$-value, without updating the prior, would lead to this configuration. Iterating rule $\mathcal{R}_1$ $n$ times, i.e. using the rule $(\mathcal{R}_1, \mathcal{R}_1, \ldots, \mathcal{R}_1)$ of (10), would of course lead to other configurations, since each newly allocated observation would provide additional information which is utilized.

13

However, one might speculate whether the fixed sample size n allocation rule $\mathcal{R}_n$ would use the (DT)-configuration. This question can be settled with the method of Lagrangian multipliers. For $k = 2$ populations it turns out to be true: a fact which has been shown already in Gupta and Miescke (1993) with a different method. For $k = 3$ , on the other hand, it can be seen that the (DT)-configuration is not a solution for $\mathcal{R}_n$ to the problem at hand. In conclusion, we can see that dealing with $k \geq 3$ populations is much more challenging than treating only $k = 2$ populations, a contrast which is typical for ranking and selection problems.

## References

[1] Abughalous, M. M. and Miescke, K. J. (1989). On selecting the largest success probability under unequal sample sizes. *Journal of Statistical Planning and Inference*, **21**, 53–68.

[2] Bechhofer, R. E. (1969). Optimal allocation of observations when comparing several treatments with a control. *Multivariate Analysis II*, P. R. Krishnaiah ed., Academic Press, New York, 465-473.

[3] Bechhofer, R. E., Hayter, A. J. and Tamhane, A. C. (1991). Designing experiments for selecting the largest normal mean when the variances are known and unequal: Optimal sample size allocation. *Journal of Statistical Planning and Inference*, **28**, 271–289.

[4] Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*. Second Edition. Springer Verlag, New York.

[5] Berger, J. O. and Deely, J. (1988). A Bayesian approach to ranking and selection of related means with alternatives to AOV methodology. *Journal of the American Statistical Association*, **83**, 364–373.

[6] Dudewicz, E.J. and Dalal, S.R. (1975). Allocation of observations in ranking and selection with unequal variances. *Sankhyā*, **B-37**, 28–78.

[7] Fong, D.K.H. and Berger, J.O. (1993). Ranking, estimation and hypothesis testing in unbalanced two-way additive models - A Bayesian approach. *Statistics and Decisions*, **11**, 1–24.

[8] Gupta, S. S. and Miescke, K. J. (1993). Bayesian look ahead one stage sampling allocations for selecting the largest normal mean. *Statistical Papers*, to appear.

[9] Gupta, S. S. and Miescke, K. J. (1988). On the problem of finding the largest normal mean under heteroscedasticity. In: *Statistical Decision Theory and Related Topics IV*, S. S. Gupta and J. O. Berger eds., Springer Verlag, New York, Vol. 2, 37–49.

[10] Gupta, S. S. and Panchapakesan, S. (1991). Sequential ranking and selection procedures. In: *Handbook of Sequential Analysis*, B. K. Ghosh and P. K. Sen eds., M. Dekker, New York, 363–379.

[11] Gupta, S. S. and Panchapakesan, S. (1979). *Multiple Decision Procedures: Theory and Methodology of Selecting and Ranking Populations*. J. Wiley, New York.

[12] Miescke, K. J. (1979). Bayesian subset selection for additive and linear loss functions. *Communications in Statistics*, A-8, 1205–1226.

[13] Whitehead, J. (1991). Sequential methods in clinical trials. In: *Handbook of Sequential Analysis*, B. K. Ghosh and P. K. Sen eds., M. Dekker, New York, 593–611.