# ON THE PROBLEM OF THE MEANS OF WEIGHTED NORMAL POPULATIONS

M. S. NIKULIN * and V. G. VOINOV‡

*An analitical problem, which arises in the statistical problem of comparing the means of two normal distributions, the variances of which as their ratio are unknown, is well-known in the mathematical statistics as the Behrens-Fisher problem. One generalization of the Behrens-Fisher problem and different aspects concerning the estimation of the common mean of several independent normal distributions with different variances are considered and one solution is proposed.*

## 1. THE BEHRENS-FISHER PROBLEM

Let $X_{ij}$   $(i = 1, \ldots, k; \; j = 1, \ldots, n_i)$ be mutually independent normally distributed random variables

$$EX_{ij} = a_i \quad \text{and} \quad \text{Var} X_{ij} = \sigma_i^2$$

In this model we have a minimal sufficient statistic

$$\mathbf{T}_{2k} = (X_1, \ldots, X_k, S_1^2, \ldots, S_k^2)^\mathsf{T},$$

*M. S. Nikulin. Laboratory of Statistical Methods. Steklov Mathematical Institute. St.Petersburg, Russia. UFR MI25. Université Bordeaux 2, France.

‡V. G. Voinov. Laboratory of Statistics. The Institute of Theoretical and Applied Mathematics. Alma-Ata, 480082. Kazakhstan.

the components of which

$$X_i = \frac{1}{n}\sum_{j=1}^{n_i} X_{ij} \quad \text{and} \quad S_i^2 = \sum_{j=1}^{n_i} (X_{ij} - X_i)^2, \quad (i = 1,\ldots,k)$$

are mutually independent statistics such that,

$$X_i \quad \text{being normal} \quad N\left(a_i, \frac{\sigma_i^2}{n_i}\right) \quad \text{and} \quad \frac{S_i^2}{\sigma_i^2} = \chi_{f_i}^2, \quad f_i = n_i - 1.$$

Let $s_i^2 = \dfrac{1}{f_i} S_i^2$, $i = 1,\ldots,k$. Since

$$E\{X_i\} = a_i \quad \text{and} \quad E\{s_i^2\} = \sigma_i^2 \quad (i = 1,\ldots,k)$$

and since the sufficient statistic $\mathbf{T}_{2k}$ is complete, we can affirm that the statistic $\mathbf{U} = (X_1,\ldots,X_k,s_1^2,\ldots,s_k^2)$ is the best (minimum variance) unbiased estimator (MVUE) for the parameter $(a_1,\ldots,a_k,\sigma_1^2,\ldots,\sigma_k^2)^{\mathsf{T}}$.

The well-known Behrens-Fisher problem consists in testing the hypothesis

$$H_0 : a_1 = a_2 = \cdots = a_k = a$$

or in constructing a confidence interval for the unknown common mean $a$ under the assumption that the solution should be written in term of the minimal sufficient statistic $\mathbf{T}_{2k}$. We note that even for $k = 2$ the exact solution of this problem has not been obtained as yet.

**Remark 1.** The statistical problem of comparing the mathematical exectations $a_1$ and $a_2$ of two normal distributions, the variances of which $\sigma_1^2, \sigma_2^2$ and their ratio $\sigma_1^2/\sigma_2^2$ are unknown, was posed by Behrens (1929). The modern formulation of this problem is due to Sir R. Fisher and is based on the concept of sufficiency, since a sufficient statistic contains the same information about the unknown parameters $a_1, \sigma_1^2, a_2, \sigma_2^2$ as the initial data $X_{11},\ldots,X_{1n_1}$ and $X_{21},\ldots,X_{2n_2}$,

$$EX_{ij} = a_i \quad \text{and} \quad \mathrm{Var}\, X_{ij} = \sigma_i^2 \quad (i = 1,2; \quad j = 1,\ldots,n_i; \; n_i \geq 2),$$

and hence only the sufficient statistic $\mathbf{T}_4 = (X_1, X_2, S_1^2, S_2^2)^{\mathsf{T}}$ needs be considered in testing hypotheses about the values of unknown parameters $a_1, \sigma_1^2, a_2, \sigma_2^2$. In particular, this approach was used by Linnik, Sudakov and Romanovsky (1964) to test the hypothesis

$$H_\delta : a_1 - a_2 = \delta,$$

where $\delta$ is a given number. In this situation the Behrens-Fisher problem reduces to construct a critical set $K_\alpha \subset \mathfrak{X}_3 \subset \mathbb{R}^3$ in the sample space $\mathfrak{X}_3$ of the three-dimentional (under the hypothesis $H_\delta$) sufficient statistic

$$\mathbf{V}_3 = (X_1 - X_2, S_1^2, S_2^2)^{\mathsf{T}}$$

such that the probability $\mathcal{P}\{\mathbf{V}_3 \in K_\alpha | H_\delta\}$ does not depend on the unknown parameters $a_i, \sigma_i^2$ and the unknown ratio $\sigma_1^2/\sigma_2^2$ and

$$\mathcal{P}\{\mathbf{V}_3 \in K_\alpha | H_\delta\} = \alpha, \quad (0 < \alpha < 0.5)$$

The question of the existence of such $K_\alpha$ was discussed during many years in the statistical literature. In 1964 Linnik, Romanovsky and Sudakov have shown (see e.g. Linnik (1966), (1968)) that if $n_1$ and $n_2$ are of different parities, then a set $K_\alpha$ to the Behrens-Fisher problem exists. But if $n_1$ and $n_2$ are of the same parities, the existence of a solution to the Behrens-Fisher problem remains an open question even today. There are many others modifications and generalizations of the Behrens-Fisher problem. For example, A. Wald posed the problem of existence of a critical set $K_\alpha$ in the sample space $\mathfrak{X}_2 \subset \mathbb{R}^2$ of the two-dimentional statistic

$$\mathbf{V}_2 = ((X_1 - X_2)/(S_1^2, S_1^2/S_2^2)^\tau.$$

The solution of this problem has not been obtained as yet. However, it is effectively possible to find a set $K_\alpha^* \subset \mathfrak{X}_2$ such that

$$\mathcal{P}\{\mathbf{V}_2 \in K_\alpha^* | H_\delta\} \cong \alpha,$$

but in this case the probability $\mathcal{P}\{\mathbf{V}_2 \in K_\alpha^* | H_\delta\}$ depends on the unknown ratio $\sigma_1^2/\sigma_2^2$. This approach is used very often in statistical applications for the practical construction of tests for the comparison of $a_1$ and $a_2$. But the statistics of these tests are not expressed in terms of sufficient statistic $\mathbf{T}_4$ and hence usually these tests are less powerful than test based on the solution of the Behrens-Fisher problems and its generalizations.

**Remark 2.** We note here that many papers have been devoted to the problem of testing the hypothesis $H_0 : a_1 = a_2 = a$, $(\delta = 0)$. For example, Welch (1947), James (1956, 1959), Brown and Forsythe (1974), Linnik (1966), Dijkstra and Werter (1981) and many others analysed the different approximations and some of them gave tables of critical values of the appropriate test statistics. Another approaches are also known. Gans (1981), for example, investigated the use of first a test for equality of the variances $\sigma_1^2 = \sigma_2^2$ and then Student's t-statistic or the alternate test, described in the above-mentioned papers, if equality is rejected.

## 2. ESTIMATION OF THE UNKNOWN COMMON MEAN

Evidently, in the case where the hypothesis $H_0$ is accepted, the problem of obtaining the best estimator for the unknown common mean arises. Naturally, we may use

both interval and point estimators of the parameter $a$. The problem of approximate interval estimator of $a$ has been considered by James (1956), Pagurova and Gursky (1979), Marič and Graybill (1979), Khatri and Shah (1981) and others. It is interesting to note the main result of Marič and Graybill (1979), which consists in setting in some sense the "exact" confidence interval for $a$ with confidence coefficient equal to any preselected value $1 - \alpha$.

The problem of point estimation of the common mean $a$ is also not simple. This is apparently due to the incompleteness of the minimal sufficient statistic. Usually the parameter $a$ is estimated by the weighted mean

(1)
$$\hat{a} = \sum_{i=1}^{k} w_i X_i,$$

where

$$W_i = \frac{\dfrac{1}{Y_i}}{\displaystyle\sum_{i=1}^{k} \dfrac{1}{Y_i}} \quad \text{and} \quad Y_i = \frac{1}{n_i(n_i - 1)} \sum_{j=1}^{n_i} (X_{ij} - X_i)^2 = \frac{S_i^2}{n_i(n_i - 1)} = \frac{1}{n_i} s_i^2,$$

which is an **unbiased** but **nonefficient estimator** of $a$, see Hinckley (1979). We note that $n_i Y_i = S_i^2/f_i = s_i^2$ is the **MVUE for** $\sigma_i^2$, $i = 1, \ldots, k$. Because of this, many authors search for more efficient estimators of $a$, see Zacks (1966), Hinckely (1979), Bhattacharya (1980) and others. If the estimator (1) is accepted, we want also to obtain its variance and the unbiased estimator of this variance. These problems have been considered by, among others, Meier (1953), Cochran and Carroll (1953), Levi and Mantel (1974), Bement and Willams (1969), Norwood and Hinkelmann (1977). Norwood and Hinkelmann showed the non-efficient estimator (1) is in some sense optimal. More specifically they showed that the variance of $\hat{a}$ is less than any of the variance $\sigma_i^2/n_i$ $(i = 1, \ldots, k)$ if and only if either $n_i > 9$ $(i = 1, \ldots, k)$ or, for some $i$, $n_i = 9$ and $n_j > 17$ $(j = 1, \ldots, k; \ j \neq i)$.

In view of the great practical importance of the problem of combining of estimators in the normal case and because of the artificiality of the estimator (1) the following approach seems interesting.

Let $X_1, \ldots, X_k, s_1^2, \ldots, s_k^2$ be mutually independent statistics, where

$$X_i \quad \text{is normal} \quad N\left(a, \frac{\sigma_i^2}{n_i}\right) \quad \text{and} \quad s_i^2 = \frac{\sigma_i^2}{f_i} \chi_{f_i}^2,$$

$f_i = n_i - 1$, $i = 1, \ldots, k$, the values $f_1, \ldots, f_k$ being fixed and parameters $a, \sigma_1^2, \ldots, \sigma_k^2$ being all unknown.

Consider the class of functions $A_k(\cdot), B_k(\cdot), C_k(\cdot), \quad k \in \mathbb{N}$, mapping $\mathbb{R}^k \times \mathbb{R}^k_+ \times \mathbb{R}^k_+$ onto $\mathbb{R}^1$, where $\mathbb{R}^k$ is the Euclidean space of $k$ dimensions, and

$$\mathbb{R}^k_+ = \{x = (x_1, \ldots, x_k) \in \mathbb{R}^k, \quad x_1 > 0, \ldots, x_k > 0\}.$$

That is, for any matrix of order $3 \times k$

(2)
$$\left\| \begin{array}{cccc} x_1 & x_2 & \cdots & x_k \\ y_1 & y_2 & \cdots & y_k \\ z_1 & z_2 & \cdots & z_k \end{array} \right\|$$

with elements satisfying the following conditions

$$|x_i| < \infty, \quad y_i > 0, \quad z_i > 0, \quad i = 1, 2, \ldots, k,$$

functions $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ give real numbers

$$A_k = A_k \left( \left\| \begin{array}{cccc} x_1 & x_2 & \cdots & x_k \\ y_1 & y_2 & \cdots & y_k \\ z_1 & z_2 & \cdots & z_k \end{array} \right\| \right), \quad B_k = B_k \left( \left\| \begin{array}{cccc} x_1 & x_2 & \cdots & x_k \\ y_1 & y_2 & \cdots & y_k \\ z_1 & z_2 & \cdots & z_k \end{array} \right\| \right),$$

$$C_k = C_k \left( \left\| \begin{array}{cccc} x_1 & x_2 & \cdots & x_k \\ y_1 & y_2 & \cdots & y_k \\ z_1 & z_2 & \cdots & z_k \end{array} \right\| \right),$$

from $\mathbb{R}^1$. Suppose, further, that functions that $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ satisfy the following conditions.

(1) $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ are symmetric with respect to rearrangement of the columns of (2).

(2) If $k = 1$, then,

$$A_1 = A_1 \left( \left\| \begin{array}{c} x_1 \\ y_1 \\ z_1 \end{array} \right\| \right) = x_1, B_1 = B_1 \left( \left\| \begin{array}{c} x_1 \\ y_1 \\ z_1 \end{array} \right\| \right) = y_1, C_1 = C_1 \left( \left\| \begin{array}{c} x_1 \\ y_1 \\ z_1 \end{array} \right\| \right) = z_1,$$

from where one can see the sense of these functions $A_k, B_k, C_k$.

(3) For any real numbers $\alpha$ and $\beta$

$$A_k \left( \left\| \begin{array}{cccc} \alpha x_1 + \beta & \alpha x_2 + \beta & \cdots & \alpha x_k + \beta \\ \alpha^2 y_1 & \alpha^2 y_2 & \cdots & \alpha^2 y_k \\ z_1 & z_2 & \cdots & z_k \end{array} \right\| \right) = \alpha A_k \left( \left\| \begin{array}{ccc} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{array} \right\| \right) + \beta,$$

97

$$B_k\left(\left\|\begin{matrix} \alpha x_1+\beta & \alpha x_2+\beta & \cdots & \alpha x_k+\beta \\ \alpha^2 y_1 & \alpha^2 y_2 & \cdots & \alpha^2 y_k \\ z_1 & z_2 & \cdots & z_k \end{matrix}\right\|\right) = \alpha^2 B_k\left(\left\|\begin{matrix} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{matrix}\right\|\right),$$

$$C_k\left(\left\|\begin{matrix} \alpha x_1+\beta & \alpha x_2+\beta & \cdots & \alpha x_k+\beta \\ \alpha^2 y_1 & \alpha^2 y_2 & \cdots & \alpha^2 y_k \\ z_1 & z_2 & \cdots & z_k \end{matrix}\right\|\right) = C_k\left(\left\|\begin{matrix} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{matrix}\right\|\right).$$

(4) For $k = n+m$ and for

$$A_n^* = A_n\left(\left\|\begin{matrix} x_1 & \cdots & x_n \\ y_1 & \cdots & y_n \\ z_1 & \cdots & z_n \end{matrix}\right\|\right), A_m^{**} = A_m\left(\left\|\begin{matrix} x_{n+1} & \cdots & x_{n+m} \\ y_{n+1} & \cdots & y_{n+m} \\ z_{n+1} & \cdots & z_{n+m} \end{matrix}\right\|\right),$$

$$B_n^* = B_n\left(\left\|\begin{matrix} x_1 & \cdots & x_n \\ y_1 & \cdots & y_n \\ z_1 & \cdots & z_n \end{matrix}\right\|\right), B_m^{**} = B_m\left(\left\|\begin{matrix} x_{n+1} & \cdots & x_{n+m} \\ y_{n+1} & \cdots & y_{n+m} \\ z_{n+1} & \cdots & z_{n+m} \end{matrix}\right\|\right),$$

$$C_n^* = C_n\left(\left\|\begin{matrix} x_1 & \cdots & x_n \\ y_1 & \cdots & y_n \\ z_1 & \cdots & z_n \end{matrix}\right\|\right), C_m^{**} = C_m\left(\left\|\begin{matrix} x_{n+1} & \cdots & x_{n+m} \\ y_{n+1} & \cdots & y_{n+m} \\ z_{n+1} & \cdots & z_{n+m} \end{matrix}\right\|\right),$$

the following relations hold:

$$A_k = A_2\left(\left\|\begin{matrix} A_n^* & A_m^{**} \\ B_n^* & B_m^{**} \\ C_n^* & C_m^{**} \end{matrix}\right\|\right), B_k = B_2\left(\left\|\begin{matrix} A_n^* & A_m^{**} \\ B_n^* & B_m^{**} \\ C_n^* & C_m^{**} \end{matrix}\right\|\right), C_k = C_2\left(\left\|\begin{matrix} A_n^* & A_m^{**} \\ B_n^* & B_m^{**} \\ C_n^* & C_m^{**} \end{matrix}\right\|\right)$$

The proposed class is not empty, since there are functions $A_{(\cdot)}, B_k(\cdot), C_k(\cdot)$ of matrices (2) satisfying conditions (1)-(4). Indeed, define functions $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ as follows

(3)
$$\left\|\begin{matrix} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{matrix}\right\| \xrightarrow{A_k(\cdot)} A_k = \dfrac{\displaystyle\sum_{i=1}^{k} \frac{x_i}{y_i}}{\displaystyle\sum_{i=1}^{k} \frac{1}{y_i}},$$

$$\left\|\begin{matrix} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{matrix}\right\| \xrightarrow{B_k(\cdot)} B_k = \dfrac{1}{\displaystyle\sum_{i=1}^{k} \frac{1}{y_i}},$$

$$\left\|\begin{matrix} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{matrix}\right\| \xrightarrow{C_k(\cdot)} C_k = \sum_{i=1}^{k} z_i.$$

Evidently, these functions satisfy condition (1). If $k = 1$, then $A_1 = x_1, B_1 = y_1, C_1 = z_1$. Hence, condtion (2) is also satisfied. Condition (3) holds, since

$$\frac{\displaystyle\sum_{i=1}^{k} \frac{\alpha x_i + \beta}{\alpha^2 y_i}}{\displaystyle\sum_{i=1}^{k} \frac{1}{\alpha^2 y_i}} = \alpha \cdot \frac{\displaystyle\sum_{i=1}^{k} \frac{x_i}{y_i}}{\displaystyle\sum_{i=1}^{k} \frac{1}{y_i}} + \beta,$$

$$\frac{1}{\displaystyle\sum_{i=1}^{k} \frac{1}{\alpha^2 y_i}} = \alpha^2 \cdot \frac{1}{\displaystyle\sum_{i=1}^{k} \frac{1}{y_i}}$$

Finally, condition (4) holds since

$$\frac{\dfrac{\displaystyle\sum_{i=1}^{n} \frac{x_i}{y_i}}{\displaystyle\sum_{i=1}^{n} \frac{1}{y_i}} }{\dfrac{1}{\displaystyle\sum_{i=1}^{n} \frac{1}{y_i}}} + \frac{\dfrac{\displaystyle\sum_{i=n+1}^{n+m} \frac{x_i}{y_i}}{\displaystyle\sum_{i=n+1}^{n+m} \frac{1}{y_i}}}{\dfrac{1}{\displaystyle\sum_{i=n+1}^{n+m} \frac{1}{y_i}}} = \frac{\displaystyle\sum_{i=1}^{k=n+m} \frac{x_i}{y_i}}{\displaystyle\sum_{i=1}^{k=n+m} \frac{1}{y_i}},$$

$$\frac{1}{\dfrac{1}{\displaystyle\sum_{i=1}^{n} \frac{1}{y_i}} + \dfrac{1}{\displaystyle\sum_{i=n+1}^{n+m} \frac{1}{y_i}}} = \frac{1}{\displaystyle\sum_{i+1}^{k=n+m} \frac{1}{y_i}},$$

$$\sum_{i=1}^{n} z_i + \sum_{i=n+1}^{n+m} z_i = \sum_{i=n+1}^{k=n+m} z_i.$$

From this it follows that functions $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ defined by (3) satisfy the conditions (1)-(4). Note that if $\alpha = -1$ and $\beta = 0$,

$$A_k\left(\left\| \begin{array}{ccc} -x_1 & \cdots & -x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{array} \right\|\right) = -A_k\left(\left\| \begin{array}{ccc} x_1 & \cdots & x_k \\ y_1 & \cdots & y_k \\ z_1 & \cdots & z_k \end{array} \right\|\right).$$

This follows from condition (3).

For convenience we denote here $\mathbf{V}_i = s_i^2 \quad (i = 1, \ldots, k)$ and let consider the statistic

$$\hat{a} = A_k \left( \left\| \begin{matrix} X_1 & \cdots & X_k \\ \mathbf{V}_1 & \cdots & \mathbf{V}_k \\ f_1 & \cdots & f_k \end{matrix} \right\| \right),$$

which by virtue of (3) may be rewritten as

$$\hat{a} = a + A_k \left( \left\| \begin{matrix} \sigma_1 \mathbf{U}_1/\sqrt{n_1} & \sigma_2 \mathbf{U}_2/\sqrt{n_2} & \cdots & \sigma_k \mathbf{U}_k/\sqrt{n_k} \\ \mathbf{V}_1 & \mathbf{V}_2 & \cdots & \mathbf{V}_k \\ f_1 & f_2 & \cdots & f_k \end{matrix} \right\| \right) = a + \xi$$

where

$$\mathbf{U}_i = \sqrt{n_i}\, \frac{X_i - a}{\sigma_i}, \quad i = 1, \ldots, k,$$

are independent standard normal variables. We now find the expected value $\mathrm{E}\xi$ of the random variable $\xi$. Under the assumption that

$$\mathbf{V}_i = v_1, \mathbf{V}_2 = v_2, \ldots, V_k = v_k, \quad (\mathbf{V}_i = s_i^2)$$

the conditional expectation $\mathrm{E}\{\xi | \mathbf{V}_1 = v_1, \mathbf{V}_2 = v_2, \ldots, \mathbf{V}_k = v_k\}$ is equal to

$$\frac{1}{(2\pi)^{k/2}} \int_{-\infty}^{+\infty} \cdots \int_{-\infty}^{+\infty} A_k \left( \left\| \begin{matrix} \sigma_1 u_1/\sqrt{n_1} & \cdots & \sigma_k u_k/\sqrt{n_k} \\ v_1 & \cdots & v_k \\ f_1 & \cdots & f_k \end{matrix} \right\| \right) e^{-\frac{1}{2}\sum_{i=1}^{k} u_i^2} \, du_1 \ldots du_k =$$

$$\frac{(-1)^k}{(2\pi)^{k/2}} \int_{-\infty}^{+\infty} \cdots \int_{-\infty}^{+\infty} A_k \left( \left\| \begin{matrix} -\sigma_1 t_1/\sqrt{n_1} & \cdots & -\sigma_k t_k/\sqrt{n_k} \\ v_1 & \cdots & v_k \\ f_1 & \cdots & f_k \end{matrix} \right\| \right) e^{-\frac{1}{2}\sum_{i=1}^{k} t_i^2} \, dt_1 \ldots dt_k =$$

$$(-1)^{2k+1} \mathrm{E}\{\xi | \mathbf{V}_1 = v_1, \mathbf{V}_2 = v_2, \ldots, \mathbf{V}_k = v_k\} =$$

$$-\mathrm{E}\{\xi | \mathbf{V}_1 = v_1, \mathbf{V}_2 = v_2, \ldots, \mathbf{V}_k = v_k\}.$$

Hence, if the expectation $\mathrm{E}\{\xi | \mathbf{V}_1 = v_1, \mathbf{V}_2 = v_2, \ldots, \mathbf{V}_k = v_k\}$ exists, it is identically equal to zero. From this fact we may conclude that $\mathrm{E}\xi = 0$ and $\mathrm{E}\hat{a} = a$. That is,

$$\mathrm{E}A_k \left( \left\| \begin{matrix} X_1 & \cdots & X_k \\ s_1^2 & \cdots & s_k^2 \\ f_1 & \cdots & f_k \end{matrix} \right\| \right) = a.$$

We may set up here the following problems.

1. Define the class of all functions $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ satisfying conditions (1)-(4) and find in this class an estimator a which has the minimal variance. Possibly, there are no other functions except $A_k(\cdot), B_k(\cdot), C_k(\cdot)$ defined by (3). These functions in such a case give the solution of the problem.

100

2. Find the variance $\mathbf{D}\hat{a}$ and its unbiased estimator.

   Now we return to our problem of the estimation of the common mean of several independent normal distributions with different variances. In particular we shall obtain the variance of (1) and its **MVUE**.

## 3. THE VARIANCE OF THE WEIGHTED MEAN AND ITS MVUE

Using the fact that

$$\sum_{i=1}^{k} w_i = 1 \quad \text{and} \quad \mathrm{E}X_i^2 = a^2 + \sigma_i^2/n_i,$$

the variance of the weighted mean (1) may be written as

$$(4) \qquad \mathbf{D}\hat{a} = \mathrm{E}\hat{a}^2 - \{\mathrm{E}\hat{a}\}^2 = \sum_{i=1}^{k} \mathrm{E}\{w_i^2\} \frac{\sigma_i^2}{n_i}.$$

Since the probability density function of the random variable $\mathbf{T}_i = 1/Y_i$ is

$$f(t_i) = \frac{b_i^{\frac{n_i-1}{2}}}{\Gamma\left(\frac{n_i-1}{2}\right)} \cdot t_i^{-\frac{n_i+1}{2}} \cdot e^{-b_i/t_i}, \quad t_i \geq 0,$$

where

$$b_i = \frac{n_i f_i}{2\sigma_i^2},$$

we have

$$\mathrm{E}\{w_j^2\} = \left\{ \prod_{i=1}^{k} \frac{b_i^{\frac{n_i-1}{2}}}{\Gamma\left(\frac{n_i-1}{2}\right)} \right\} \int_0^{\infty} \cdots \int_0^{\infty} \frac{t_j^2}{\left(\sum_{i=1}^{k} t_i\right)^2} \cdot$$

$$(5) \qquad \prod_{i=1}^{k} t_i^{-\frac{n_i+1}{2}} \exp\left\{ -\sum_{i=1}^{k} \frac{b_i}{t_i} \right\} dt_1 \ldots dt_k.$$

It is well known that

$$\left(\sum_{i=1}^{k} t_i\right)^{-2} = \int_0^\infty u \exp\left\{-u\sum_{i=1}^{k} t_i\right\} du.$$

Substituting this expression into (5) and rearranging the order of integration, we get

$$E\{w_j^2\} = \left\{\prod_{i=1}^{k} \frac{b_i^{\frac{n_i-1}{2}}}{\Gamma\left(\frac{n_i-1}{2}\right)}\right\} \int_0^\infty u\,du \int_0^\infty \cdots \int_0^\infty t_j^{-\frac{n_j-3}{2}} e^{-u t_j\,\frac{b_j}{t_j}} \cdot$$

$$\cdot \prod_{i\neq j}\left(t_i^{-\frac{n_j+1}{2}} e^{-u t_i\,\frac{b_i}{t_i}}\right) dt_1 \ldots dt_k =$$

$$= b_j 2^{\frac{n-3k}{2}-1}\left\{\prod_{i=1}^{k} \frac{b_i^{\frac{n_i-1}{4}}}{\Gamma\left(\frac{n_i-1}{2}\right)}\right\} \int_0^\infty x^{\frac{n-k}{2}+1} \mathcal{K}_{\frac{n_i-5}{2}}\left(x\sqrt{b_i}\right) \cdot$$

$$\cdot \prod_{i\neq j} \mathcal{K}_{\frac{n_i-1}{2}}\left(x\sqrt{b_i}\right) dx,$$

where $n = n_1 + n_2 + \cdots + n_k$ and $\mathcal{K}_i(\cdot)$ is the modified Bessel function. Using this result and formula (4), one may express the variance $\mathbf{D}\hat{a}$ of the weighted mean $\hat{a}$ as

$$\mathbf{D}\hat{a} = \frac{1}{2^{\frac{n-3k}{2}+1}}\left\{\prod_{i=1}^{k} \frac{b_i^{\frac{n_i-1}{4}}}{\Gamma\left(\frac{n_i-1}{2}\right)}\right\} \sum_{j=1}^{k}\left(\frac{n_j-1}{2}\right) \int_0^\infty x^{\frac{n-k}{2}+1} \cdot$$

(6)
$$\cdot \mathcal{K}_{\frac{n_j-5}{2}}\left(x\sqrt{b_j}\right) \cdot \prod_{i\neq j} \mathcal{K}_{\frac{n_i-1}{2}}\left(x\sqrt{b_j}\right) dx.$$

102

If $k = 2$, for example, see Gradshteyn and Ryzhik (1980),

$$\mathbf{D}\hat{a} = \frac{b_1^{\frac{n_1-1}{4}} b_2^{\frac{n_1-1}{4}}}{2^{\frac{n-4}{2}} \Gamma\left(\frac{n_1-1}{2}\right) \Gamma\left(\frac{n_2-1}{2}\right)} \left\{ \frac{b_1\sigma_1^2}{n_1} \int_0^\infty x^{\frac{n}{2}} \mathcal{K}_{\frac{n_1-5}{2}}\left(x\sqrt{b_1}\right) \cdot \right.$$

$$\left. \cdot \mathcal{K}_{\frac{n_2-1}{2}}(x\sqrt{b_2}) \, dx + \frac{b_2\sigma_2^2}{n_2} \int_0^\infty x^{\frac{n}{2}} \mathcal{K}_{\frac{n_1-1}{2}}\left(x\sqrt{b_1}\right) \mathcal{K}_{\frac{n_2-5}{2}}\left(x\sqrt{b_2}\right) dx \right\} =$$

$$= \frac{b_2^{\frac{n_1-1}{2}} \Gamma\left(\frac{n-2}{2}\right) \Gamma\left(\frac{n_2+3}{2}\right) \sigma_1^2}{b_1^{\frac{n_2-1}{2}} \Gamma\left(\frac{n+2}{2}\right) \Gamma\left(\frac{n_2-1}{2}\right) n_1} \, {}_2F_1\left(\frac{n-2}{2} \cdot \frac{n_2+3}{2}; \frac{n+2}{2}; 1 - \frac{b_2}{b_1}\right) +$$

$$+ \frac{b_2^{\frac{n_1-1}{2}} \Gamma\left(\frac{n-2}{2}\right) \Gamma\left(\frac{n_1+3}{2}\right) \sigma_2^2}{b_1^{\frac{n_2-1}{2}} \Gamma\left(\frac{n+2}{2}\right) \Gamma\left(\frac{n_1-1}{2}\right) n_2} \, {}_2F_1\left(\frac{n-2}{2}, \frac{n_2-1}{2}; \frac{n+2}{2}; 1 - \frac{b_2}{b_1}\right).$$

Here ${}_mF_1(\alpha, \beta; \gamma; x)$ is the Gauss hypergeometric function. One may easily verify that this expression is identical to one given by Nair (1980) who used another approach.

Now we may find an unbiased estimate of (6). This problem reduces to constructing a statistic the expected value of which is exactly equal to $\mathbf{D}\hat{a}$. Since $\mathbf{D}\hat{a}$ does not depend on $a$, the complete sufficient statistic of our problem is the vector $(Z_1, Z_2, \ldots, Z_k)^{\mathsf{T}}$, where $Z_i = f_i/2Y_i$ $(i = 1, 2, \ldots, k)$.

The probability density function of this random variable $Z_i$ is

$$f(z_i) = \frac{[(n_i-1)b_i]^{\frac{n_i-1}{2}}}{2^{\frac{n_i-1}{2}} \Gamma\left(\frac{n_i-1}{2}\right)} z_i^{-\frac{n_i+1}{2}} \exp\left\{-\frac{(n_i-1)b_i}{2z_i}\right\}, \quad z_i \geq 0, \quad i = 1, \ldots, k.$$

Using tables of Gradshteyn and Ryzhik (1980), it is easy to verify that the unbiased estimates of

$$(7) \qquad b_i^{\frac{n_i-1}{4}} \mathcal{K}_{\frac{n_i-1}{2}}\left(x\sqrt{b_i}\right) \quad \text{and} \quad \frac{n_i-1}{2} b_i^{\frac{n_i-1}{2}} \mathcal{K}_{\frac{n_i-5}{2}}\left(x\sqrt{b_i}\right)$$

are

$$\frac{2^{\frac{n_i-3}{2}} \Gamma\left(\frac{n_i-1}{2}\right)}{x^{\frac{n_i-1}{2}}} \exp\left\{-\frac{x^2 z_i}{2(n_i-1)}\right\}$$

and

$$(8) \qquad \frac{2^{\frac{n_i-7}{4}} \Gamma\left(\frac{n_i-1}{2}\right)}{x(n_1-1)^{\frac{n_i-7}{4}}} Z_i^{\frac{n_i-3}{4}} W_{-\frac{n_i-3}{4},\frac{n_i-5}{4}}\left(\frac{x^2 Z_i}{2(n_i-1)}\right) e^{-\frac{x^2 Z_i}{4(n_i-1)}}$$

respectively, where $W_{\mu,\nu}(\cdot)$ is the Whittaker function. Substituting in (6) unbiased estimates (8) instead of (7) and performing the integration with respect to $x$, we get (see, also, Nikulin and Voinov, (1983), Voinov and Nikulin (1993))

$$\mathbf{D}\hat{a} = \Sigma_{j=1}^{k} \frac{Z_j^{\frac{n_j-3}{2}}}{2(n_j-1)^{\frac{n_j-3}{4}}} \left(\sum_{i=1}^{k} \frac{Z_i}{n_i-1}\right)^{-\frac{n_j-3}{2}} {}_2F_1\left(\frac{n_j-1}{2}, \frac{n_j-3}{2}; \frac{n_j+1}{2}; \frac{\Sigma_{i\neq j}\frac{z_i}{n_i-1}}{\Sigma_{i=1}^{k}\frac{z_i}{n_i-1}}\right) =$$

$$= \Sigma_{j=1}^{k} \frac{1}{Y_j^{\frac{n_j-3}{2}}} \left(\sum_{i=1}^{k} \frac{1}{Y_i^2}\right)^{-\frac{n_j-1}{2}} {}_2F_1\left(\frac{n_j-1}{2}, \frac{n_j-3}{2}; \frac{n_j+1}{2}; \frac{\Sigma_{i\neq j}\frac{1}{Y_i^2}}{\Sigma_{i=1}^{k}\frac{1}{Y_i^2}}\right).$$

Since

$${}_2F_1(a,b;c;x) = (1-x)^{c-a-b} {}_2F_1(c-a,c-b;c;x),$$

the unique minimum variance unbiased estimate of the variance of the weighted mean may be written as

$$(9) \qquad \mathbf{D}\hat{a} = \sum_{j=1}^{k} \frac{c^2}{Y_j^2} {}_2F_1\left(1,2;\frac{n_j+1}{2};1-\frac{c}{Y_j^2}\right), \quad \text{where} \quad c = \frac{1}{\Sigma_{i=1}^{k}\frac{1}{Y_i^2}}.$$

If
$$\min_{1 \le i \le k} n_i \longrightarrow \infty, \quad \text{then} \quad {}_2F_1\left(1,2;\frac{n_i+1}{2};x\right) \longrightarrow 1$$

and expression (9) converges to $\left(\sum_{i=1}^{k} n_i/\sigma_i^2\right)^{-1}$ which is the variance of the best linear unbiased estimate of the weighted mean, if $\sigma_i^2$ are known. One may easily verify that expression (9) to, within terms of order $O(1/n_i^2)$, coincides with the approximately unbiased estimate of Meier (1953).

## ACKNOWLEDGEMENTS

## REFERENCES

[1] **Behrens W.U.** (1929). *Landwirtsch. Jahresber.* **68, 6**, 807–837.

[2] **Bement, T.R.** and **Williams, J.S.** (1969). "Variance of weighted regression estimators when sampling errors are independent and heteroscedastic". *JASA*, **64**, 1369–1382.

[3] **Bhattacharya, C.G.** (1980). "Estimation of a common mean and recovery of interblock information". *Ann. Statisti.*, **8, 1**, 205–211.

[4] **Brown, M.B.** and **Forsythe, A.B.** (1974). "The small sample behavior of some statistics which test the equality of several means". *Technometrics*, **16**, 129–132.

[5] **Cochran, W.G.** and **Carroll, S.P.** (1953). "A sampling investigation of the efficiency of weighting inversely as the estimated variance". *Biometrics*, **9**, 447–459.

[6] **Dijkstra, J.B.** and **Werter, P.S.P.J.** 1(1981). "Testing the equality of several means when the population variances are unequal". *Commun. in. Statist - Simula. Computa*, **B 10, 6**, 557–569.

[7] **Gans, D.J.** (1981). "Use of a preliminary test in comparing two sample means". *Commun. in Statist. - Simula. Computa*, **B 10, 2**, 163–174.

[8] **Gradshteyn, I.S.** and **Ryzhik, I.M.** (1980). "Table of integrals, series, and products". *Academic Press*, New York.

[9] **Hinckley, D.V.** (1979). "A note on the weighted means problem". *Scand. J. Statist.*, **6, 1**, 37–40.

[10]  **James, G.S.** (1959). "The BehrensFisher distribution and weighted means". *Journal Royal Statist. Soc*, **B 21**, 73–90.

[11]  **James, G.S.** (1956). "On the accuracy of weighted means and ratios". *Biometrics*, **43**, 304–321.

[12]  **Khatri, C.G.** and **Shah, K.R.** (1981). "Interval estimation of the common mean". *Commun. in Statist. - Simula. Computa.*, **B 10, 2**, 99–107.

[13]  **Levy, P.S.** and **Mantel, N.** (1974). "Combining unbiased estimates - A further examination of some old estimators". *J. Statist. Comput. Simul.*, **3**, 147–160.

[14]  **Linnik, Yu.V., Romanovsky, I.V., Sudakov, V.N.** (1964). "A norandomized homogeneous test in the Behrens-Fisher problem". *Dokl. Acad. Nauk SSSR*, **155, 6**, 1262–1264.

[15]  **Linnik, Yu.V.** (1966). "Randomized homogeneous tests for the Behrens-Fisher problem". In: *Selected Trans.in Math.Stat. and Prob.*, **6**, 207–217.

[16]  **Linnik, Yu.V.** (1968). "Statistical problems with nuisance parameters". *Amer. Math. Soc.*, (translated from Russian).

[17]  **Marič, N.** and **Graybill, F.A.** (1979). "Small sample confidence intervals on the common mean of two normal distributions with unequal variances". *Commun. in Statist. - Theor. Math.*, **A 8, 13**, 1255–1269.

[18]  **Meir, P.** (1953). "Variance of a weighted mean". *Biometrics*, **9, 1**, 59–73.

[19]  **Nair, K.A.** (1980). "Variance and distribution of the Graibill-Deal estimator of the common mean of two normal population". *Ann. of Statist.*, **8, 1**, 212–216.

[20]  **Nikulin, M.S.** and **Voinov, V.G.** (1983). "On the problem of the weighted mean of several normal populations". *Preprint LOMI AN SSSR*, E183, Leningrad.

[21]  **Norwood, T.E.** and **Hinkelmann, K**. (1977). "Estimating the common mean of several normal populations". *Ann. of Statist.*, **5, 5**, 1047–1050.

[22]  **Pagurova, V.I.** (1968). "On a comparison of means of two normal samples". *Theory of probability and its applications*, **13, 3**, 527–534.

[23]  **Pagurova, V.I.** and **Gursky, V.V.** (1979). "A confidence interval for the common mean of several normal distributions". *Theory of Probability and its applications*, **24, 4**, 882–888.

[24]  **Voinov, V.G.** and **Nikulin, M.S.** (1993). "Unbiased estimators and their applications". **1**. Univariate case. *Kluwer*, Holland.

[25]  **Welch, B.L.** (1947). "The generalization of "Student's" problem when several different population variances are involved". *Biometrika*, **34**, 28–35.

[26]  **Zacks, S.** (1966). "Unbiased estimation of the common mean of two normal distributions based on small samples of equal size". *JASA*, **61**, 467–476.