# Note on Selection from a Multivariate Normal Population

## By A. C. AITKEN.

The problem of statistical "selection" is concerned with the alteration induced in a frequency distribution in several variables by an alteration of the parameters in a subsection of the distribution. It may be illustrated by a simple trivariate case, as follows:

From a population characterised by variables $x$, $y$, $z$, correlated and normally distributed, with means 0, 0, 0, variances $\sigma_1^2$, $\sigma_2^2$, $\sigma_3^2$ and product variances $r_{12}\sigma_1\sigma_2$, $r_{13}\sigma_1\sigma_3$, $r_{23}\sigma_2\sigma_3$, a sub-population is extracted by selection in $x$ alone, in such a way that after selection $x$ is still normally distributed, but with mean $h$ and variance $s^2$. It is required to determine the new values, in the selected population, of the means and variances of $y$ and $z$, and of the product variances.

This type of problem was first considered (1) by K. Pearson, and was solved by him for the general case of $n$ normally correlated variables of which $p$ variables are selected, that is, are subjected to an alteration of their $p$ means, $p$ variances and $\frac{1}{2}p(p-1)$ product variances, the remaining $n-p$ variables undergoing an induced alteration in respect of their corresponding parameters. It was shown (2) by H. E. Soper that such problems, which include as special cases problems of partial correlation and regression, could be treated conveniently by the use of moment generating functions. This method, which is the same in principle as the Heaviside operational calculus and the use of Fourier, Hankel and other integral transforms in pure analysis, proceeds by setting up a duality between a frequency function and its exponential transform, defined by, *e.g.*

$$G(a, \beta, \gamma) = \int\!\!\!\int\!\!\!\int_{-\infty}^{\infty} \phi(x, y, z)\, e^{ax+\beta y+\gamma z}\, dx\, dy\, dz,$$

where $\phi(x, y, z)\, dx\, dy\, dz$ is the frequency differential. For example in the case of three normally correlated variables, with variances standardized to unity, we have the very convenient function

$$G(a, \beta, \gamma) = \exp\{\tfrac{1}{2}(a^2 + \beta^2 + \gamma^2 + 2r_{12}\,a\beta + 2r_{13}\,a\gamma + 2r_{23}\,\beta\gamma)\}.$$

Since $a$ interprets $x$, Soper proceeds by gathering all the terms in $a$ into a squared expression, thus,

$$\exp \{ \tfrac{1}{2} (a + r_{12} \beta + r_{13} \gamma)^2 + \tfrac{1}{2} \beta^2 (1 - r_{12}^2) + \beta \gamma (r_{23} - r_{12} r_{13}) + \tfrac{1}{2} \gamma^2 (1 - r_{13}^2) \}.$$

He then replaces the exponential factor involving $a$ by its equivalent integral involving frequency in $x$ alone, thus

$$G(a, \beta, \gamma) = \left\{ (2\pi)^{-\frac{1}{2}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}x^2} e^{(a + r_{12}\beta + r_{13}\gamma)x} dx \right\} \exp \{ \tfrac{1}{2}\beta^2 (1 - r_{12}^2) + \dots \}.$$

Evidently the effect of selection in $x$, producing a new mean $h$ and a new variance $s^2$, will be to alter the frequency differential from

$$(2\pi)^{-\frac{1}{2}} \exp(-\tfrac{1}{2}x^2)\, dx \quad \text{to} \quad (2\pi)^{-\frac{1}{2}} \exp\{ -\tfrac{1}{2}(x - h)^2/s^2 \}\, dx/s.$$

Making this substitution and integrating back, Soper obtains at once the new moment generating function

$$\exp \{ h(a + r_{12}\beta + r_{13}\gamma) + \tfrac{1}{2}[s^2 a^2 + (1 - r_{12}^2 \overline{1 - s^2})\beta^2 + \dots]$$
$$+ r_{12} s^2 a\beta + r_{13} s^2 a\gamma + (r_{23} - r_{12} r_{13} \overline{1 - s^2})\beta\gamma \}.$$

By mere inspection the coefficients of $a$, $\beta$, $\gamma$ in the exponent give the new means as $h$, $r_{12} h$, $r_{13} h$; the coefficients of $a^2/2!$, $\beta^2/2!$, $\gamma^2/2!$ give the new variances as

$$s^2, \quad 1 - r_{12}^2 (1 - s^2), \quad 1 - r_{13}^2 (1 - s^2),$$

and the coefficients of $a\beta$, $a\gamma$, $\beta\gamma$ give the new product variances as

$$r_{12} s^2, \quad r_{13} s^2, \quad r_{23} - r_{12} r_{13} (1 - s^2).$$

These results furnish the complete solution of this trivariate problem; and the case $s = 0$ gives the solution of the corresponding problem in partial correlation.

### General Solution of the Multivariate Problem.

By a slight simplification of the above approach, and the free use of the notation and methods of matrices, a perspicuous formulation of the solution of Pearson's general problem may be obtained.

The $n$ variables may be viewed synoptically[1] as a vector $x$,

---

[1] Vectors $x$ and $a$ are matrices of one column of elements, $x'$ and $a'$ are matrices of one row obtained from $x$ and $a$ by "transposition." Quadratic forms are written $x'Ax$ or $a'Aa$, where $A$ is symmetric, i.e. $A' = A$.

interpreted in the moment generating function by a vector $a$. It may be assumed without loss of generality that the vector of means is the null vector, and that the $n$ variances have been standardized to have unit value. The moment generating function before selection is therefore, in a suitable notation,

$$G\,(a) = (2\pi)^{-\frac{1}{2}n}\,|\,R\,|^{-\frac{1}{2}} \int_{-\infty}^{\infty} \exp\left(-\tfrac{1}{2}x'\,R^{-1}\,x\right) e^{a'x}\,dx = \exp\left(\tfrac{1}{2}a'\,Ra\right), \quad (1)$$

where $R$ is the matrix $[r_{ij}]$ of total correlation coefficients, and where the integral is $n$-ple, $dx$ standing for $dx_1\,dx_2\,\ldots\,.\,dx_n$.

*Lemma.*  Since

$$\int_{-\infty}^{\infty} \phi\,(x - h)\, e^{a'x}\,dx = e^{a'h} \int_{-\infty}^{\infty} \phi\,(x)\, e^{a'x}\,dx, \qquad (2)$$

a vector change of origin from the vector $0$ to the vector $h$ induces a factor $\exp\,(a'h)$ in the moment generating function.

In preparation for selection according to the first $p$ of the $n$ variables, the matrix $R$ is partitioned thus,

$$R = \begin{bmatrix} R_{pp} & R_{pq} \\ \hline R_{qp} & R_{qq} \end{bmatrix}, \qquad (3)$$

that is, into submatrices of orders $p \times p$, $p \times q$, $q \times p$, $q \times q$, where $p + q = n$. The vector $a$ is correspondingly partitioned, its first $p$ elements being regarded as a vector which without confusion can still be denoted by $a$, provided the remaining $q$ elements are distinguished as a vector $\beta$. The quadratic form $\tfrac{1}{2}a'Ra$ thus appears in partitioned shape as

$$\begin{aligned} &\tfrac{1}{2}\,(a'R_{pp}\,a + 2a'R_{pq}\,\beta + \beta'R_{qq}\,\beta) \\ &= \tfrac{1}{2}\{(a + R_{pp}^{-1}\,R_{pq}\,\beta)'R_{pp}\,(a + R_{pp}^{-1}\,R_{pq}\,\beta) + \beta'(R_{qq} - R_{qp}\,R_{pp}^{-1}\,R_{pq})\,\beta\}, \end{aligned} \quad (4)$$

the vector $a$ being segregated in a first quadratic term.

The exponential of this first quadratic term may be replaced by a multiple integral, involving a frequency differential (with matrix $R_{pp}^{-1}$ in the integrand; *cf.* (1) above) in the first $p$ variables only, multiplied by $\exp\{(a + R_{pp}^{-1}\,R_{pq}\,\beta)'x\}$.

Now let selection operate, (i) by altering the vector of means of these first $p$ variables from $0$ to $h$. By the Lemma noted above in (2), the moment generating function acquires a factor

$$\exp\{(a + R_{pp}^{-1}\,R_{pq}\,\beta)'\,h\}.$$

This yields the first general result, that all the new means are given as coefficients of the $a_i$ and $\beta_i$ in the elements of the vector

$$(a + R_{pp}^{-1} R_{pq} \beta)'\, h. \tag{5}$$

Next, (ii) let the variances and product variances of the first $p$ variables be altered. The effect of this is that $R_{pp}^{-1}$, the matrix in the integrand of frequency for the first $p$ variables, is replaced by some other symmetric matrix, say $V_{pp}^{-1}$. On integrating back (*cf.* again (1) above) the *reciprocal* matrix $V_{pp}$ must therefore replace $R_{pp}$ as matrix of the first quadratic term in (4). Hence the whole quadratic form (4) is altered by selection to

$$\tfrac{1}{2}\{(a + R_{pp}^{-1} R_{pq} \beta)'\, V_{pp}\, (a + R_{pp}^{-1} R_{pq} \beta) + \beta'\, (R_{qq} - R_{qp} R_{pp}^{-1} R_{pq})\, \beta\}$$

$$= \tfrac{1}{2} \{a'\, V_{pp}\, a + 2a'\, V_{pp}\, R_{pp}^{-1} R_{pq} \beta + \beta'(R_{qq} - R_{qp} R_{pp}^{-1} R_{pq} + R_{qp} R_{pp}^{-1} V_{pp} R_{pp}^{-1} R_{pq})\, \beta\}. \tag{6}$$

The results of Pearson and Soper are contained in the two statements (5) and (6). The second may be expressed thus: the effect of selection is to change

$$\begin{bmatrix} R_{pp} & R_{pq} \\ \hline R_{qp} & R_{qq} \end{bmatrix} \text{ into } \begin{bmatrix} V_{pp} & V_{pp} R_{pp}^{-1} R_{pq} \\ \hline R_{qp} R_{pp}^{-1} V_{pp} & R_{qq} - R_{qp} (R_{pp}^{-1} - R_{pp}^{-1} V_{pp} R_{pp}^{-1}) R_{pq} \end{bmatrix}. \tag{7}$$

*Example.* To deduce the formulae for selection in one of three variables.

$$\text{Here} \qquad R = \begin{bmatrix} 1 & r_{12} & r_{13} \\ \hline r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{bmatrix}, \qquad V_{pp} = [s_1^2].$$

By the general result, the matrix after selection becomes

$$\begin{bmatrix} [s_1^2] & [s_1^2][1][r_{12} \quad r_{13}] \\ \hline \text{etc.} & \begin{bmatrix} 1 & r_{23} \\ r_{23} & 1 \end{bmatrix} - \begin{bmatrix} r_{12} \\ r_{13} \end{bmatrix} [1][r_{12} \quad r_{13}] + \begin{bmatrix} r_{12} \\ r_{13} \end{bmatrix} [s_1^2][r_{12} \quad r_{13}] \end{bmatrix}$$

$$= \begin{bmatrix} s_1^2 & r_{12} s_1^2 & r_{13} s_1^2 \\ r_{12} s_1^2 & 1 - r_{12}^2 (1 - s_1^2) & r_{23} - r_{12} r_{13} (1 - s_1^2) \\ r_{13} s_1^2 & r_{23} - r_{12} r_{13} (1 - s_1^2) & 1 - r_{13}^2 (1 - s_1^2) \end{bmatrix}$$

and here all the results of Soper's example are visible.  As for the means, we have

$$(a + R_{pp}^{-1} R_{pq} \beta)'h \equiv (a_1 + r_{12} a_2 + r_{13} a_3) h,$$

which gives them as $h$, $r_{12} h$, $r_{13} h$, again as found earlier.

————

## REFERENCES.

1.   K. Pearson, *Phil. Trans. Roy. Soc.*, London, **200** A (1902), 1-66.
2.   H. E. Soper, *Frequency Arrays* (Cambridge, 1922), 20-21.