# A NOTE ON THE DISTRIBUTION OF THE STUDENTISED D²-STATISTIC

By

S. N. ROY

STATISTICAL LABORATORY

## INTRODUCTION

In an earlier paper[1] published by the author jointly with Mr. R. C. Bose the Studentised D²-Statistic was defined, its statistical use and object were fully set forth and its sampling distribution was worked out. It is the object of the present note to considerably simplify the derivation of the distribution obtained there, by a twist in the geometrical arguments employed in that paper and a consequent change of procedure. This will incidentally throw fresh light on the use of hyperspace geometry in tackling this class of problems. To make this note intelligible even without reference to the earlier paper[1], or in other words, to make this note practically independent of that, we have to take over three small sections from that paper, where definitions and notations were laid down and certain preliminaries were also introduced which will suffice for our present purpose.

## §1. PRELIMINARIES.

Consider two samples $\Sigma$ and $\Sigma'$ of sizes $n$ and $n'$ from two multivariate normal populations $\Pi$ and $\Pi'$ with the same set of variances and covariances $a_{ij}$ $(i, j = 1, 2, \ldots \ldots \ldots p)$ where $a_{ij} = \rho_{ij} \cdot \sigma_i \cdot \sigma_j$, $\sigma_i$ and $\sigma_j$ being the standard deviations for the $i$-th and $j$-th characters respectively, and $\rho_{ij}$ the correlation co-efficient between the $i$-th and $j$-th characters.[1] The matrix $\| a_{ij} \|$ will be said to be the common dispersion matrix for the two populations. Let $a_{ij}$, $a'_{ij}$ $(i, j = 1, 2, \ldots \ldots p)$ denote the respective variances and covariances of the samples $\Sigma$ and $\Sigma'$, so that $\| a_{ij} \|$ and $\| a'_{ij} \|$ are their respective dispersion matrices. Let $\alpha_i$, $\alpha'_i$ $(i = 1, 2, \ldots \ldots p)$ be the means for the $i$-th character for the populations $\Pi$ and $\Pi'$ and let $a_i$, $a'_i$ denote the corresponding quantities for the samples $\Sigma$ and $\Sigma'$.

Let us set

$$c_{ij} = \frac{n a_{ij} + n' a'_{ij}}{n + n'} \qquad \ldots \quad (1 \cdot 1)$$

Let $c^{ij}$ as usual denote the minor of $c_{ij}$ in the determinant $| c_{ij} |$ divided by the determinant itself. A like definition holds for $a^{ij}$. Then the Studentised D²-statistic is defined by

$$p \, D^2 = c^{11}(a_1 - a'_1)^2 + c^{22}(a_2 - a'_2)^2 + \ldots \ldots c^{pp}(a_p - a'_p)^2 + 2c^{12}(a_1 - a'_1)(a_2 - a'_2) + \ldots \ldots \ldots 2c^{p-1,p}(a_{p-1} - a'_{p-1})(a_p - a'_p) \quad \ldots \quad (1 \cdot 15)$$

Likewise if $\Delta^2$ is the population value of $D^2$, then

$$p\,\Delta^2 = a^{11}(\alpha_1 - \alpha'_1)^2 + a^{22}(\alpha_2 - \alpha_2)^2 + \ldots\ldots a^{pp}(\alpha_p - \alpha'_p)^2$$
$$+ 2a^{12}(\alpha_1 - \alpha'_1)\,(\alpha_2 - \alpha'_2) + \ldots\ldots 2a^{p-1,p}(\alpha_{p-1} - \alpha'_{p-1})\,(\alpha_p - \alpha'_p) \qquad \ldots \quad (1\cdot2)$$

We know that the joint distribution of the sample readings $x_{ik}$, $x'_{ik'}(i = 1, 2, \ldots \ldots \ldots p;\ k = 1, 2, \ldots \ldots n;\ k' = 1, 2, \ldots \ldots n')$ is

$$\text{const. } e^{-\frac{1}{2}\sum\limits_{i=1}^{p}\sum\limits_{i=1}^{p} a^{ij}\{n(a_i - \alpha_i)\,(a_j - \alpha_j) + n'(a'_i - \alpha'_i)\,(a'_j - \alpha'_j) + (nd_{ij} + n'a'_{ij})\}}$$
$$\times \Pi\,dx_{ik}\,\Pi\,dx'_{ik'} \quad \ldots \quad (1\cdot25)$$

where $\Pi\,dx_{ik}$ stands for $dx_{11}\,dx_{12}\,dx_{21}\ldots\ldots dx_{pn}$ and a similar meaning attaches to $\Pi dx'_{ik'}$.

Sample $\Sigma$ can be represented in the usual Fisherian space $S_n$ of $n$ dimensions, by the points with co-ordinates

$$(x_{i1}, x_{i2}, \ldots \ldots x_{in}),\ i = 1, 2, \ldots \ldots p. \qquad \ldots \quad (1\cdot3)$$

or what is the same thing, by $p$ vectors $x_i$ joining the points to the origin. We take another space $S'_{n_1}$ of $n'$ dimensions absolutely orthogonal to the former space, and represent in it the sample $\Sigma'$ by $p$ other similar vectors. Let

$$y_{ik} = x_{ik} - a_{i.},\ y'_{ik'} = x'_{ik'} - a'_{i.}, \qquad \ldots \quad (1\cdot35)$$

where $i = 1, 2, \ldots \ldots p;\ k = 1, 2, \ldots \ldots n;\ k' = 1, 2, \ldots \ldots n'$

Let $y_i$, $y'_i$ denote the vectors, with components $(y_{i1}, y_{i2}, \ldots \ldots y_{in})$ and $(y'_{i1}, y'_{i2}, \ldots \ldots y'_{in'})$ lying in the space $S_n$ and $S'_n$ respectively. Then the vectors $y_i$ lie in a flat $S_{n-1}$ of $n-1$ dimensions, perpendicular to the equiangular line in $S_n$. Similar considerations apply to the vector $y'_i$.

Let O be the origin of co-ordinates and let $M_i$ be the point on the equiangular line in $S_n$ such that $OM_i = \dfrac{1}{\sqrt{n}}$ times the projection of $x_i$ on the equiangular line. Then $O\,M_i = a_i$. In the same way we can find $M'_i$ on the other equiangular line such that $O\,M'_i = a'_i$. Also if $y_i$. $y_j$ is the scalar product of the vectors $y_i$ and $y_j$, then clearly

$$a_{ij} = (y_i.\ y_j)/n,\quad a'_{ij} = (y'_i.\ y'_j)/n' \qquad \ldots \quad (1\cdot4)$$

Let us now take a new set of $(n+n')p$ variable $a_i$, $a'_i$, $z_{ik}$, $z'_{ik'}$ ($i = 1, 2, \ldots \ldots p;\ k = 1, 2, \ldots \ldots n-1;\ k' = 1, 2, \ldots \ldots n'-1$)

such that

$$(i)\quad a_i = \frac{1}{n}\sum_{k=1}^{n} x_{ik}\ (i = 1, 2, \ldots \ldots \ldots \ldots \ldots \ldots p)$$

$$(ii)\quad a'_i = \frac{1}{n'}\sum_{k'=1}^{n'} x'_{ik'}\ (i = 1, 2, \ldots \ldots \ldots \ldots \ldots \ldots p)$$

(iii) $z_{ik}$ $(k = 1, 2, \ldots \ldots \ldots n-1)$ are the components of $y_i$ along any $n-1$ mutually orthogonal lines in $S_{n-1}$. Then $z_{ik}$ is naturally a linear function of

$$x_{i1}, x_{i2}, \ldots \ldots \ldots \ldots \ldots \ldots x_{in}$$

(iv) Similar considerations apply to $z'_{ik'}$.

The distribution (1·25) can now be written in the form

$$\text{con t} \times e^{-\frac{1}{2} \sum\limits_{i=1}^{p} \sum\limits_{j=1}^{p} a^{ij} \{n(a_i - a_i)(a_j - a_j) + n'(a'_i - a'_i)(a'_j - a'_j) + (n \, a_{ij} + n' \, a'_{ij})\}}$$

$$\times \prod_{i=1}^{p} da_i \prod_{i=1}^{p} da'_i \, \Pi \, dz_{ik} \, \Pi dz'_{ik'} \quad \ldots \quad (1·45)$$

It should be noted that $a_{ij}$'s are expressible purely in terms of $z_{ik}$'s and $a'_{ij}$'s are expressible purely in terms of $z'_{ik}$'s

Next we introduce the new variables $a_i - a'_i$ and $a_i + a'_i$ in place of $a_i$ and $a'_i$ and integrating for $a_i + a'_i$ we get the distribution in the form

$$\text{const} \times e^{-\frac{1}{2} \sum\limits_{i=1}^{p} \sum\limits_{j=1}^{p} a^{ij} \left[ \frac{n}{2} \{(a_i - a'_i) - (a_1 - a'_1)\} |(a_j - a'_j) - (a_2 - a'_2)\} + N c_{ij} \right]}$$

$$\times \prod_{i=1}^{p} d(a_i - a'_i) \, \Pi \, dz_{ik} \, \Pi \, dz'_{ik'} \quad \ldots \quad (1·5)$$

where $c_{ij}$ is given by (1·1) and

$$\frac{2}{\bar{n}} = \frac{1}{n} + \frac{1}{n'}, \quad N = n + n' \quad \ldots \quad (1·55)$$

In the plane of the equiangular lines of $S_n$ and $S'_{n'}$ let $R_i$ be the point, whose projections on the equiangular lines coincide with $M_i$ and $M'_i$. Then if $Q_i$ is the projection of $R_i$ on $O Y$, the external bisector of the equiangular lines

$$O Q_i = \frac{1}{\sqrt{2}} (a_i - a'_i) \quad (i = 1, 2, \ldots \ldots \ldots p) \quad \ldots \quad (1·6)$$

Let the vector $v_i$ be the resultant of the vectors $y_i$ and $y'_i$. Then it is easily seen that

$$N \, c_{ij} = v_i . v_j \quad \ldots \quad (1·65)$$

where the dot denotes the scalar product.

We may note that the spaces $S_{n-1}$, $S'_{n'-1}$ containing the vectors $y_i$, $y'_i$ respectively $(i = 1, 2, \ldots \ldots p)$ are orthogonal to one another, as also to $O Y$. Hence the new vectors $v_i$ are also orthogonal to $O Y$, and lie in the space $S_{n+n'-2}$ which comprises both $S_{n-1}$ and $S'_{n'-1}$.
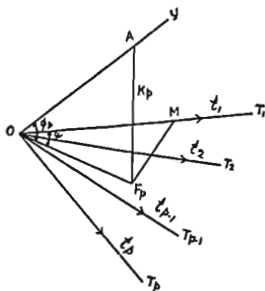
Let $t_i$ be the resultant of the vectors $v_i$ and the vector $\bar{n}^{\frac{1}{2}} Q Q_i$, where $\bar{n}$ is given by (1·55).

Then

$$t_i \cdot t_i = \frac{\bar{n}}{2} (a_i - a'_i) (a_i - a'_j) + (na_{ij} + n'a_{ij}')$$

$$= \frac{\bar{n}}{2} (a_i - a'_i) (a_i - a'_j) + N c_{ij} = g_{ij} \text{ (say)} \qquad (1·7)$$

The distribution (1·5) now takes the form

$$\text{const} \times e^{-\frac{1}{2} p \Delta^2 - \frac{1}{2} \sum_{i=1}^{p} \sum_{j=1}^{p} a^{ij} \{ g_{ij} - \bar{n} (a_i - a'_i) (a_j - a'_j) \}}$$

$$\times \Pi \ d(a_i - a'_i) \ \Pi \ dz_{ik} \ \Pi \ dz'_k. \qquad (1·75)$$



Let $T_i$ denote the extremity of the vector $t_i$.

Let OA be the unit vector along OY. We shall denote it by $i$. Let $F_i$ denote the foot of the perpendicular from A to the space $(T_1 \ T_2 \ldots\ldots\ldots T_i)$. Let the length $AF_i$ be denoted by $K_i$

Let us now consider the length $K_p$. We shall show that it is very closely connected with $D^2$

$$\text{Now } K^2_p = \frac{\text{Vol } (i, t_1, t_2,\ldots\ldots\ldots t_p)}{\text{Vol } (t_1, t_2,\ldots\ldots\ldots t_p)} \qquad \ldots \ (1·8)$$

But Vol $(i, t_1, t_2,\ldots\ldots\ldots t_p)$ is the same as the volume formed by the unit vector QA, and the projections of the vectors $t_1, t_2,\ldots\ldots\ldots t_p$ on the space perpendicular to OA, which we have earlier called $S_{n,n-2}$. But these projections, from the way in which

$\ell_1, \ell_2,\cdots\ldots\ell_p$ have been derived, are at once seen to be $v_1, v_2,\ldots\ldots\ldots v_p$. Hence the numerator in (1·95) is numerically equal to vol $(v_1, v_2,\ldots\ldots\ldots v_p)$

Therefore, $K^2_p = \dfrac{|\ v_1 \cdot\ v_1\ |}{|\ \ell_1 \cdot\ \ell_1\ |}$

$= \dfrac{|\ N\ c_{ij}\ |}{|\ N\ c_{ij} + \dfrac{\bar{n}}{2}\ (a_i - a'_i)\ (a_j - a'_j)\ |}$   from (1·65) and (1·7)

$= \dfrac{|\ N\ c_{ij}\ |}{|\ N\ c_{ij}\ | + \dfrac{\bar{n}}{2}\ N^{p-1}\ \Sigma\ C_{ij}(a_i - a'_i)\ (a_j - a'_j)}$   ... (1·81)

$C_{ij}$ denoting the minor of $c_{ij}$ in $|\ c_{ij}\ |$

Therefore from (1·15) $K^2_p = \dfrac{1}{1 + \dfrac{p\bar{n}}{2N}\ D^2}$   ... ... (1·82)

Denoting by $\varphi_p$ the angle $AOF_p$, i.e. the angle between the unit vector $i$, and the flat $(\ell_1, \ell_2,\ldots\ldots\ldots\ell_p)$ we have, since $k_p = \sin \varphi_p$

$$D^2 = \dfrac{2N}{p\ \bar{n}}\ \cot^2 \varphi_p \qquad \ldots (1·83)$$

If a statistic similar to $D^2$ were constructed for the first $i$ variates only, then we could denote it by $D^2_i$; we should then have had

$$K^2_i = \dfrac{1}{1 + \dfrac{i\bar{n}}{2N}\ D^2_i} \qquad \ldots (1·84)$$

and

$$D^2_i = \dfrac{2N}{i\bar{n}}\ \cot^2 \varphi_i \qquad \ldots (1·85)$$

It should be noted that consistently with this notation $D^2$ should be preferably replaced by $D^2_p$. Let M be the foot of the perpendicular from $F_p$ to $\ell_1$ and let $v$ be the angle $F_p$ OM.

## §2.  DERIVATION OF THE SAMPLING DISTRIBUTION OF $D^2$

It was shown in section 2 of the earlier paper[1] that $D^2$ and $\Delta^2$ defined in the present note by equations (1·15), (1·2), and $\sum\limits_{i=1}^{p}\ \sum\limits_{j=1}^{p}\ a^{ij}\ a_{ij}$, $\sum\limits_{i=1}^{p}\ \sum\limits_{j=1}^{p}\ a^{ij}\ (a_i - a_j)\ (a_i - a_j)$ occurring in the present note in the relation (1·75) are all invariant under any linear transformation. It was further shown in section 8 of the earlier paper[1] that a special linear transformation could be constructed such that

$$\left.\begin{array}{l} a_{ij} = 0,\ \text{if}\ i \neq j,\ a_{ii} = 1 \\[4pt] (a_1 - a'_1) = p\ \Delta^2 \\[4pt] (a_1 - a_i) \approx 0\ \text{for}\ i = 2,\ 3,\ldots\ldots\ldots\ldots\ldots p \end{array}\right\} \qquad \ldots (2·1)$$

Therefore, $a^{ij} = 0$ if $i \neq j$, $a^{ii} = 1$   ... (2·11)

We may consider $\Delta$ to be given by the positive root of $\Delta^2$ and $\sqrt{p}$ the positive root of $p$. Then by suitably naming the populations as first and second we can write

$$a_1 - a'_1 = \Delta \sqrt{p} \qquad \qquad \ldots \ (2\cdot12)$$

Therefore without any loss of generality we can take the distribution in $(1\cdot75)$ as

$$\text{const} \times e^{-\frac{1}{2}\bar{n}\,p\Delta^2 - \frac{1}{2}\sum\limits_{i=1}^{p} g_{ii} - \frac{\bar{n}}{2}\sqrt{p}\cdot\Delta(a_1 - a'_1)} \qquad \Pi d(a_i - a'_i)\ \Pi dz_{ik}\ \Pi dz'_{ik}. \qquad \ldots \ (2\cdot13)$$

The density factor in $(2\cdot13)$ does not contain any mean difference except the first. Therefore integrating out for all mean differences $(a_i - a'_i)$, $(i = 2, 3, \ldots\ldots p)$ we have the distribution in the form

$$\text{const} \times e^{-\frac{1}{2}\bar{n}\,p\Delta^2 - \frac{1}{2}\sum\limits_{i=1}^{p} g_{ii} - \frac{\bar{n}}{2}\sqrt{p}\,\Delta(a_1 - a'_1)} \qquad d(a_1 - a'_1)\Pi dz_{ik}\ \Pi dz'_{ik}. \qquad \ldots \ (2\cdot14)$$

Now it is clear that the volume element $d(a_1 - a'_1)\ \Pi dz_{ik}\ \Pi dz'_{ik}$. is really the volume described by the points $(A, T_1, T_2, \ldots\ldots\ldots T_p)$ in a space of $(n + n' - 1)$ dimensions, when the points are given whatever freedom we please but subject to the restriction that the variates will have to lie within infinitesimal ranges round about the values occurring in the density factor in $(2\cdot14)$. Let us denote the lengths $OT_1, OT_2, \ldots\ldots OT_p$ by $l_1, l_2, \ldots\ldots l_p$. It is now clear from the definitions of $(a_1 - a'_1)$ and $g_{ii}$ in $(1\cdot6)$ and $(1\cdot7)$ and from the figure in section 1 that

$$g_{ii} = l^2_i, \text{ and } (\sqrt{\bar{n}}/2)\ (a_1 - a'_1) = \text{projection of } OT_1 \text{ on } OY = l_1 \cos \phi_p \cos \psi \quad \ldots \ (2\cdot2)$$

The density factor in $(2\cdot14)$ now reduces to

$$\text{const} \times e^{-\frac{1}{2}\bar{n}\,p\Delta^2 - \frac{1}{2}\sum\limits_{i=1}^{p} l^2_i - (\bar{n}\,p/2)^{1/2}\ \Delta\ l_1 \cos \phi_p \cos \psi} \qquad \ldots \ (2\cdot21)$$

Give now the points $(A, T_1, T_2, \ldots\ldots\ldots T_p)$ freedom to move subject to the restrictions that $l_1$ lies between $l_1$ and $l_1 + dl_1$ $(i = 1, 2, \ldots\ldots p)$ and $\phi_p$ and $\psi$ lie respectively between $\phi_p$ and $\phi_p + d\phi_p$ and $\psi$ and $\psi + d\psi$. Let this be done in two stages. It is clear that keeping the vectors $(t_1, t_2, \ldots l_p)$ fixed the volume element described by $A$ is

$$\text{const.}\ (\sin \phi_p)^{N-p-2} (\cos \phi_p)^{p-1} (\sin \psi)^{p-2}\ d\phi_p\ d\psi \qquad \ldots \ (2\cdot22)$$

where $N$ of course is $n + n'$

Release now the points $T_1, T_2, \ldots\ldots\ldots T_p$ subject to the restrictions mentioned before; the volume element described by these points is

$$\text{const.}\ (l_1, l_2, \ldots\ldots l_p)^{N-2}\ dl_1,\ dl_2, \ldots\ldots\ldots dl_p \qquad \ldots \ (2\cdot23)$$

Therefore, the joint volume element described by $(\Lambda, T_1, T_2, \ldots\ldots\ldots T_p)$ i.

$$\text{const. } (\sin \varphi_p)^{N-p-2} (\cos \varphi_p)^{p-1} (\sin \psi)^{p-2} d\varphi_p \, d\psi \, (t_1, t_2, \ldots\ldots t_p)^{N-2} \Pi dt_1 \qquad \ldots \ (2\cdot24)$$

If we restrict now each of $\varphi_p$ and $\psi$ to lie between 0 and $\pi/2$ we easily see from (2·21) and (2·24) that the distribution (2·14) now reduces to

$$\text{const} \times e^{-\frac{1}{2}\bar{n} \, p \Delta^2 - \frac{1}{2} \sum_{i=1}^{N} t_i^2} \cosh (c \Delta t_1 \cos \varphi_p \cos \psi) (\sin \varphi_p)^{N-p-2} (\cos \varphi_p)^{p-1} (\sin \psi)^{p-2}$$

$$\times d\varphi_p \, d\psi \, (t_1, t_2, \ldots \ldots t_p)^{N-2} \prod_{i=1}^{p} dt_1 \qquad \ldots \ (2\cdot25)$$

where $c^2 = \bar{n} \, p/2$ $\qquad \ldots \ (2\cdot255)$

Integrating out for $dt_2, dt_3, \ldots\ldots\ldots dt_p$ we have

$$\text{const} \times e^{-\frac{1}{2}\bar{n} \, p \Delta^2 - \frac{1}{2} t_1^2} \cosh (c \Delta t_1 \cos \varphi_p \cos \psi) (\sin \varphi_p)^{N-p-2} (\cos \varphi_p)^{p-1} (\sin \psi)^{p-2}$$

$$\times d\varphi_p \, d\psi \, t_1^{N-2} dt_1 \qquad \ldots \ (2\cdot26)$$

It should be noted that $t_1$ varies from 0 to $\infty$ and both $\varphi_p$ and $\psi$ vary from 0 to $\pi/2$

Using now the well known relation[1]

$$I_\nu (z) = \frac{2(z/2)^\nu}{\Gamma(\nu+\frac{1}{2}) \, \Gamma(\frac{1}{2})} \int_0^{\pi/2} \cosh (z \cos \theta) \sin^{2\nu} \theta \, d\theta \qquad \ldots \ (2\cdot3)$$

which is valid for $R(\nu+\frac{1}{2})>0$, for the special case $\nu=0$, and integrating (2·26) over $\psi$ fr 0 to $\pi/2$ the distribution reduces to

$$\text{const} \times e^{-\frac{1}{2} t_1^2} (t_1)^{N-2-(p-2)/2} I_{\frac{p-2}{2}} (c \Delta t_1 \cos \varphi_p) (\sin \varphi_p)^{N-p-2} (\cos \varphi_p)^{p/2} dt_1 \, d\varphi_p \qquad \ldots \ (2\cdot35)$$

Using now another well known relation[2]

$$\int_0^\infty I_\nu (at) \, e^{-b^2 t^2} \, t^{\mu-1} dt = \frac{\Gamma(\frac{1}{2} \nu + \frac{1}{2} \mu)}{2b^\mu \, \Gamma(\nu+1)} \cdot {}_1F_1\left(\frac{1}{2} \nu + \frac{1}{2} \mu, \nu+1, \frac{a^2}{4b^2}\right) \qquad \ldots \ (2\cdot4)$$

for the special case $b^2 = \frac{1}{2}$, $\nu = (p-2)/2$, $\mu = N-1-(p-2)/2$, $a = c \, \Delta \cos \varphi_p$, which do not violate the restrictions on the validity of the formula, and integrating (2·35) over $t_1$ from 0 to $\infty$, the distribution of $\varphi_p$ is obtained in the form

$$\text{const} \times (\sin \varphi_p)^{N-p-2} (\cos \varphi_p)^{p-1} {}_1F_1\left(\frac{N-1}{2}, \frac{p}{2}, \frac{c^2 \, \Delta^2 \cos^2 \varphi_p}{2}\right) d\varphi_p \qquad \ldots \ (2\cdot45)$$

But remembering from (1·83) and (2·255) that

$$D^2 = \frac{N}{c^2} \cot^2 \varphi_p \qquad \qquad ... \quad (2·46)$$

and therefore,

$$\left.\begin{aligned} d\varphi_p &= \frac{c \sqrt{N}.\, d\, D^2}{2\,(D^2)^{\frac{1}{2}}\,(N + c^2\, D^2)} \\[2mm] \cos \varphi_p &= \frac{c(D^2)^{\frac{1}{2}}}{(N + c^2\, D^2)^{\frac{1}{2}}}, \quad \sin \varphi_p = \frac{N}{(N + c^2\, D^2)^{\frac{1}{2}}} \end{aligned}\right\} \qquad ... \quad (2·5)$$

we have the distribution of $D^2$ given in the form

$$\text{const.} \ \frac{(D^2)^{\frac{(p-3)/2}{2}}}{(N + c^2\, D^2)} \frac{N - p - 2}{2} \ {}_1F_1\left(\frac{N-1}{2}, \ \frac{p}{2}, \ \frac{c^4\,\Delta^2}{2} \ \frac{D^2}{N + c^2 D^2}\right) d(D^2) \qquad .. \quad (2·55)$$

BIBLIOGRAPHY.

1.  BOSE, R. C. AND ROY, S. N.:   The Distribution of the Studentised $D^2$-Statistic.   The Proceedings of the Indian Statistical Conference, Calcutta, 1938, pp. 19—38.

2.  WATSON, G. N.:   Theory of Bessel Functions, p. 79.

3.  Loc. Cit., p. 100 and p. 393.