

## MISCELLANEOUS NOTES

### TEST OF SIGNIFICANCE FOR INTRA-CLASS CORRELATION WHEN FAMILY SIZES ARE NOT EQUAL.

By R. P. BHARGAVA  
*Statistical Laboratory, Calcutta.*

#### INTRODUCTION.

The test of significance for the intra-class correlation when family sizes are equal was given by Fisher in his book "Statistical Methods for Research Workers".

In practice, we often meet with cases where family sizes are unequal. Thus, if we consider a certain character among brothers in a family, it often happens that the number of brothers in the different families is not the same. In such cases, some of the problems that arise are.

- (1) test of significance, that is, whether population correlation coefficient is equal to any assigned value  $\rho$
- (2) to find the fiducial limits for  $\rho$
- (3) to determine the best estimate of population correlation  $\rho$ .
- (4) to get a suitable estimate of population mean  $m$ .
- (5) to estimate intra-class correlation from family means of two unequal family sizes, the individual observations of the family being not known.

#### 2. TESTS OF SIGNIFICANCE.

Let us suppose there are  $n$  families and let the  $k_i$  individuals of the  $i$ th family give the observations

$$x_{i1}, x_{i2}, \dots, x_{ik_i} \quad (i = 1, 2, \dots, n)$$

We consider  $x_{i1}, x_{i2}, \dots, x_{ik_i}$  to be a random observation from a  $k_i$  variate normal population in which each of the variates has the same mean  $m$ , and the same standard deviation  $\sigma$ , the correlation between any two variates being equal to  $\rho$ . It may be noted that as  $m$ ,  $\sigma$  and  $\rho$  are not functions of  $i$ , each family has the same mean  $m$ , the same standard deviation  $\sigma$  and the same correlation coefficient  $\rho$ .

Let  $\bar{x}_i$  be the sample mean of the  $i$ th family. The probability of occurrence of the observation in the  $i$ th family is given by

$$\text{Const.} \cdot \frac{1}{\sigma^{k_i}(1-\rho)} \left[ \prod_{j=1}^{k_i} (x_{ij} - m)^2 - \frac{\rho}{1+k_i-1} \sum_{l, l'=1}^{k_i} (x_{il} - m)(x_{il'} - m) \right] \prod_{j=1}^{k_i} dx_{ij}$$

$$\text{where Const.} = \left( \sqrt{\frac{2\pi\sigma}{1-\rho}} \right)^{-k_i} \frac{-(k_i-1)/2}{(1-\rho)^{k_i-1} (1+k_i-1)^{k_i-1} \rho}$$

This can be written as

$$\text{Const.} \cdot \frac{1}{\sigma^{k_i}(1-\rho)} \left[ \sum_{j=1}^{k_i} (x_{ij} - \bar{x}_i)^2 + \frac{(1-\rho)}{1+k_i-1} \frac{k_i}{\rho} (\bar{x}_i - m)^2 \right] \prod_{j=1}^{k_i} dx_{ij} \quad (2.1)$$

Hence the joint distribution of the sets of observations of our sample can be written as

$$\text{Const.} \cdot \frac{1}{\sigma^{k_i}(1-\rho)} \left[ \sum_{i=1}^n (x_{i1} - \bar{x}_i)^2 + (1-\rho) \sum_{i=1}^n \frac{k_i}{1+k_i-1} \rho (\bar{x}_i - m)^2 \right] \prod dx_{ij} \quad (2.2)$$

$$\text{Let} \quad \sigma^2 = \frac{k_i(1-\rho)}{1+k_i-1} \rho, \quad \Sigma \sigma^2 \bar{x}_i = \Sigma \sigma^2 \bar{x} \quad (2.3)$$

Then (2.2) can be written as

$$\text{Const.} \cdot \frac{1}{2\sigma^2(1-\rho)} \{ \sum(x_{11} - \bar{x}_1)^2 + \sum c_i^2 (\bar{x}_i - \bar{x})^2 + \sum c_i^2 (\bar{x} - m)^2 \} \quad |d x_{11} \quad (2.4)$$

From this it is easily seen that

$$F = \frac{\sum c_i^2 (\bar{x}_i - \bar{x})^2}{(n-1)} \cdot \frac{\sum(k_i - 1)}{\sum(x_{11} - \bar{x}_1)^2} \quad (2.5)$$

is distributed in the usual F distribution with  $n_1 = n-1$  and  $n_2 = (\sum k_i - 1)$  degrees of freedom.

This F statistic can be used for testing whether our families (sizes being unequal) came from multi-variate normal populations having population correlation  $\rho$ . Thus as a particular case if we desire to know, whether our sample of  $n$  families with unequal sizes came from multivariate normal populations having  $\rho = .5$ , all that we have to do is, to calculate  $c_i$  and  $\bar{x}$  by putting  $\rho = .5$ , and then calculate the value of F by using (2.5) and consult the F tables.

Thus we can with the help of the F distribution answer any of the following questions (i) whether population correlation is significantly different from .5 (say), (ii) whether the population correlation coefficient is significantly greater than .5 and (iii) whether the population correlation coefficient is significantly less than .5.

These questions are the same as those which arise in testing the significance of intra-class correlation for families of equal size. But it is useful to note that for answering question (i) we have to consider both the tail-ends of the F distribution, while for question (ii) we have only to consider the tail-end on the right and for (iii) the tail-end on the left.

It is of interest to note that (i) Fisher's Z-test for testing the significance of intra-class correlation for equal family size and (ii) Fisher's technique of analysis of variance for testing the significance of intra-class correlation when family sizes may not be equal follow also from this test.

(i) Thus if in (2.5), we put  $k_1 = k_2 = \dots = k_k = k$  any then F becomes

$$F = \frac{1-\rho}{1+k-1\rho} \cdot \frac{\sum k(\bar{x}_i - \bar{x})^2}{n-1} \cdot \frac{n(k-1)}{\sum(x_{11} - \bar{x}_1)^2} \\ = \frac{1-\rho}{1+k-1\rho} \cdot \frac{1+k-1}{1-r} \cdot \frac{n}{n-1} \quad (2.7)$$

where  $\bar{x}$  is the grand mean and  $r$  is the sample intra-class correlation. And it is easily seen that Fisher's Z is the same as  $\frac{1}{2} \log F$ .

(ii) Again putting  $\rho = 0$ , F is equal to

$$F = \frac{\sum k_i (\bar{x}_i - \bar{x})^2}{n-1} \cdot \frac{\sum(k_i - 1)}{\sum(x_{11} - \bar{x}_1)^2}$$

which is the same as that used for testing the significance of between and within variation in analysis of variance.

### 3. FIDUCIAL LIMITS OF $\rho$

Next we consider the problem of ascertaining fiducial limits of population  $\rho$  for the case of unequal family sizes. The usual procedure for finding the lower fiducial limit is to first find (say) the 5% point of F corresponding to  $n_1 = n-1$ ,  $n_2 = \sum(k_i - 1)$ . Let us call this point as  $F_5$  then from the equality

$$\frac{\sum c_i^2 (\bar{x}_i - \bar{x})^2}{n-1} \cdot \frac{\sum(k_i - 1)}{\sum(x_{11} - \bar{x}_1)^2} = F_5$$

we solve for  $\rho$ . From this we get that  $\rho > \rho_5$  in 95% cases. This provides us with the lower limit of  $\rho$ . But in our case, the solution does not come out very easily, hence we shall follow the method of trial and error and graduation. Let us calculate values of F of (2.6) for various values of  $\rho$  say  $\rho = .1, .2, .3, .4, .5, .6, .7, .8, .9$ .

### SIGNIFICANCE FOR INTRA-CLASS CORRELATION

This can be presented in a tabular form as follows:—

$\rho$	$F(\rho)$	$F_\alpha (n_1, n_2)$
.1	$1 < F(\rho = .1)$	$> F_\alpha (n_1 = n - 1, n_2 = \Sigma(k_i - 1))$
.2	$1 < F(\rho = .2)$	$> F_\alpha (n_1 = n - 1, n_2 = \Sigma(k_i - 1))$
.3	$1 < F(\rho = .3)$	$< F_\alpha (n_1 = n - 1, n_2 = \Sigma(k_i - 1))$
.4	$1 = F(\rho = .4)$	
.5	$1 > F(\rho = .5)$	$> \frac{1}{F_\alpha} (n_1 = \Sigma k_i - 1, n_2 = n - 1)$
.6	$1 > F(\rho = .6)$	$> \frac{1}{F_\alpha} (n_1 = \Sigma k_i - 1, n_2 = n - 1)$
.7	$1 > F(\rho = .7)$	$< \frac{1}{F_\alpha} (n_1 = \Sigma k_i - 1, n_2 = n - 1)$
.8	$1 > F(\rho = .8)$	$< \frac{1}{F_\alpha} (n_1 = \Sigma k_i - 1, n_2 = n - 1)$
.9	$1 > F(\rho = .9)$	$< \frac{1}{F_\alpha} (n_1 = \Sigma k_i - 1, n_2 = n - 1)$

In the above table  $F_\alpha = (n_1, n_2)$  in column (3) indicates the 5% point of F table corresponding to  $n_1$  and  $n_2$  degrees of freedom. Column (2) gives the value of F corresponding to particular values of  $\rho$ .

In the above hypothetical table it is clear that the fiducial limits of  $\rho$  are given by

$$.2 < \rho < .7$$

This gives a rough idea, and some more numerical values of  $\rho$  between .2 and .3, and .7 and .9 will (with the help of graduation or otherwise) give closer approximation for the fiducial limits.

#### 4. BEST ESTIMATE OF $\rho$

Next comes the problem of obtaining the best estimate of  $\rho$ . The value of  $\rho$  in the above table which makes F equal to unity is the best estimate of  $\rho$ , in the sense of maximising the likelihood of  $z = \frac{1}{2} \log F$ .

Thus from the above table we see that  $\rho = .4$  gives  $F = 1$ . Thus our point estimate of  $\rho$  from the given sample comes as .4. This may be called unbiased sample intra-class correlation  $r$ .

It is clear from the way we have proceeded that our method of determining fiducial limits and point estimate of  $\rho$  coincide with the usual methods when family sizes are equal.

#### 5. ESTIMATE OF $\bar{x}$

From (4) it is clear that the distribution of  $\bar{x}$  is

$$\text{Const. } \cdot \frac{1}{2\sigma^2} \frac{k_1}{1+k_1-1} \frac{1}{\rho} (\bar{x}-m)^2 \quad d\bar{x} \quad (5.1)$$

where

$$\bar{x} = \frac{\sum \frac{k_i \bar{x}_i}{1+k_i-1}}{\sum \frac{k_i}{1+k_i-1}}$$

Thus  $\bar{x}$  is a sufficient statistic to estimate  $m$ . Thus if population  $\rho$  is known, we can calculate the best estimate of  $m$ . But, as is usual, if  $\rho$  is not known, we can estimate  $m$ , substituting for  $\rho$ , the sample intra-class correlation  $r$ , and this gives us a consistent estimate for  $m$ . Thus

$$\bar{m} = \frac{\sum \frac{k_i \bar{x}_i}{1+k_i-1} r}{\sum \frac{k_i}{1+k_i-1}}$$

is a consistent sample estimate of  $m$ . The use of this very often arises in crop surveys, budget enquiries etc

## B. ESTIMATION OF INTRA-CLASS CORRELATIONS FROM TWO SETS OF FAMILIES

Let us suppose that there are  $n$  families of size  $k$  and  $n'$  families of size  $k'$ , ( $k' < k$ ) coming respectively from  $k$  and  $k'$  variate normal populations having the same mean  $m$ , the same standard deviation  $\sigma$ , and the same correlation  $\rho$ . The individual readings in a family are not known but the family means (or totals) are known. The problem is to get an estimate of  $\rho$ .

Let  $\bar{x}_i$  be the mean of the  $i$ th family of size  $k$ ; and let  $\bar{x}'_i$  be the mean of the  $i$ th family of size  $k'$ .

Then

$$n\bar{x} = \sum_{i=1}^n x_i, \text{ and } n'\bar{x}' = \sum_{i=1}^{n'} x'_i \quad (6.1)$$

$$\text{also } (n-1)s^2 = \sum k(\bar{x}_i - \bar{x})^2, \text{ and } (n'-1)s'^2 = \sum k'(\bar{x}'_i - \bar{x}')^2 \quad (6.2)$$

Then from (2.4), it is easily seen that

$$F = \frac{1 + (k'-1)\rho}{1 + (k-1)\rho} \frac{s'^2}{s^2} \quad (6.3)$$

is distributed in the usual F distribution with  $n_1 = n-1$  and  $n_2 = n'-1$  degrees of freedom.

From this, very often we can get an estimate of  $\rho$ . The numerical value of  $\rho$ , lying between  $\frac{-1}{k-1} < \rho < 1$  which makes  $F$  of (6.3) equal to unity, is the best estimate of  $\rho$  from family means.

It is further seen that whenever, the ratio of the two mean squares  $s'^2/s^2$  is  $>k/k'$  no such numerical value of  $\rho$  lying in  $\frac{-1}{k-1} < \rho < 1$  exists which will make  $F$  of (6.3) equal to unity; but on the other hand, as is usually the case, if the ratio  $s'^2/s^2$  is  $<k/k'$  such a numerical value of  $\rho$  always exists. And in such cases this estimate is given by

$$r = -1/(k-1) + \frac{s'^2}{s^2 - s'^2} \cdot (k-k')$$

## SUMMARY.

- (1) We have given a test of significance (unknown) for testing whether our sample of  $n$  families, family size being unequal, has come from multivariate normal populations having population correlation coefficient equal to  $\rho$ .
- (2) We have given a method of deriving the fiducial limits of  $\rho$ .
- (3) We have given a new definition for intra-class correlation when family size is unequal. This definition reduces to the usual definition for the equal family size and is the best estimate of  $\rho$ .
- (4) We have derived a consistent estimate for the population mean  $m$ .
- (5) We have discussed a connected problem of estimation of intra-class correlation from family means of two unequal family sizes when individual readings of the family are not known but only family means or totals are known.

## REFERENCES.

- (1) FISHER, R. A. (1944) Statistical Methods for Research workers, 6th Edition.
- (2) ————— (1938) Statistical Theory of Estimation, Calcutta University Readership Lectures.
- (3) KENDALL, M. G. (1943) The Advanced Theory of Statistics, Vol. I.

Paper received: 4 April 1946.