

ON CERTAIN PROBABLE ERRORS AND CORRELATION COEFFICIENTS OF MULTIPLE FREQUENCY DISTRIBUTIONS WITH SKEW REGRESSION.

By L. ISSERLIS, D.Sc.

(1) In the systematic investigation of the statistical constants of multiple correlation and of their probable errors, it is important to have to hand the probable errors and the mutual correlations of the more fundamental constants—the means, the standard deviations and the correlation coefficients. For the case in which the frequency distribution follows the normal law this need is supplied in the memoir by Pearson and Filon entitled “On the Probable Errors of Frequency Constants and on the Influence of Random Selection on Variation and Correlation*.”

As regards the more general case in which the regression is skew, the probable error of a correlation coefficient was first given by Sheppard (*Phil. Trans.* Vol. 192, A, p. 128).

The probable error of a mean and the correlation between deviations in the value of the mean and that of a standard deviation, or of a correlation coefficient, and the correlation between two standard deviations, are given by Pearson (*Biometrika*, Vol. ix. 1913, pp. 1-10).

For reference we give here the results for the case of normal distributions obtained by Pearson and Filon in the memoir referred to above:

$$\begin{aligned} \Sigma_{\sigma_1} &= \sigma_1/\sqrt{2n} \dots\dots\dots(1), \\ \Sigma_{r_{12}} &= (1 - r_{12}^2)/\sqrt{n} \dots\dots\dots(2), \\ R_{\sigma_1\sigma_2} &= r_{12}^2 \dots\dots\dots(3), \\ R_{\sigma_1r_{12}} &= r_{12}/\sqrt{2} \dots\dots\dots(4), \\ R_{\sigma_1r_{23}} &= \frac{r_{12}(r_{13} - r_{12}r_{23}) + r_{13}(r_{12} - r_{12}r_{23})}{\sqrt{2} \cdot (1 - r_{23}^2)} \dots\dots\dots(5), \\ R_{r_{12}r_{23}} &= r_{23} - \frac{r_{12}r_{13}(1 - r_{23}^2 - r_{12}^2 - r_{13}^2 + 2r_{12}r_{13}r_{23})}{2(1 - r_{12}^2)(1 - r_{13}^2)} \dots\dots\dots(6), \\ R_{r_{12}r_{34}} &= \frac{\left\{ \begin{aligned} &(r_{13} - r_{12}r_{23})(r_{34} - r_{23}r_{34}) + (r_{14} - r_{34}r_{13})(r_{23} - r_{12}r_{13}) \\ &+ (r_{12} - r_{14}r_{34})(r_{34} - r_{12}r_{14}) + (r_{14} - r_{12}r_{34})(r_{23} - r_{34}r_{34}) \end{aligned} \right\}}{2(1 - r_{12}^2)(1 - r_{34}^2)} \dots\dots\dots(7) \dagger. \end{aligned}$$

* *Phil. Trans.* Vol. 191, A (1898), pp. 229-311.

† *Phil. Trans.* Vol. 191, A, Equations (xv)-(xviii), (xxxvi), (xxxvii) and (xl).

In the present paper the corresponding results are obtained for the case of skew regression. The method employed is different and by supposing the regression to be linear and the distribution to be normal, a confirmation is obtained of the above results which in Pearson and Filon's memoir depend on very complicated analysis.

(2) We may begin by discussing the correlation that exists between deviations from their means in the case of two correlation coefficients r_{xy} and r_{xt} .

We have
$$r_{xy} = p_{xy} / \sqrt{(p_x p_y)} \dots\dots\dots(8),$$

$$r_{xt} = p_{xt} / \sqrt{(p_x p_t)} \dots\dots\dots(9),$$

where $p_{x^l y^m z^n t^k}$ is employed to denote the mixed moment of orders l, m, n, k in the variables, taken about the means so that

$$\frac{dr_{xy}}{r_{xy}} = \frac{dp_{xy}}{p_{xy}} - \frac{1}{2} \frac{dp_x}{p_x} - \frac{1}{2} \frac{dp_y}{p_y} \dots\dots\dots(10),$$

and
$$\frac{dr_{xt}}{r_{xt}} = \frac{dp_{xt}}{p_{xt}} - \frac{1}{2} \frac{dp_x}{p_x} - \frac{1}{2} \frac{dp_t}{p_t} \dots\dots\dots(11).$$

It is clear that we shall require the correlations between any one of p_{xy}, p_x, p_y and any one of p_{xt}, p_x and p_t . It will suffice to find the correlation between p_{xy} and p_{xt} .

Now
$$Np_{xy} = SS \{n_{xy} (x - \bar{x}) (y - \bar{y})\} \dots\dots\dots(12),$$

$$\therefore Ndp_{xy} = SS \{dn_{xy} (x - \bar{x}) (y - \bar{y})\} + SS \{-n_{xy} (y - \bar{y}) d\bar{x}\} + SS \{-n_{xy} (x - \bar{x}) d\bar{y}\} \dots\dots(13),$$

or
$$Ndp_{xy} = SS \{dn_{xy} XY\} \dots\dots\dots(14),$$

if we denote the total population by $N, x - \bar{x}$ by $X, y - \bar{y}$ by Y and remember that

$$S \{(x - \bar{x}) n_{xy}\} = S \{(y - \bar{y}) n_{xy}\} = 0.$$

Similarly
$$Ndp_{xt} = SS \{dn_{xt} ZT\} \dots\dots\dots(15).$$

The mean value of $dn_{xy} dn_{xt}$ in many samples is the mean value of

$$\begin{aligned} & SS \{dn_{x_1 y_1 z_1 t_1}\} \times SS \{dn_{x_2 y_2 z_2 t_2}\} \\ &= -\frac{1}{N} SSSS \{dn_{x_1 y_1 z_1 t_1} dn_{x_2 y_2 z_2 t_2}\} + n_{x_1 y_1 z_1 t_1} (1 - n_{x_1 y_1 z_1 t_1} / N) \dots(16), \end{aligned}$$

where in the fourfold summation the term

$$(dn_{x_1 y_1 z_1 t_1})^2$$

is omitted.

But clearly the right-hand member of (16) reduces to

$$n_{x_1 y_1 z_1 t_1} - n_{x_1 y_1} n_{z_1 t_1} / N,$$

hence the mean value of

$$N^2 dp_{xy} dp_{xt} \text{ is } N(p_{xyzt} - p_{xy} p_{xt}) \dots\dots\dots(17).$$

Putting $t = z$ we deduce from (17)

$$\text{Mean value of } dp_{xy} dp_z = (p_{xyz} - p_{xy}p_z)/N \dots\dots\dots(18),$$

and putting $y = x$ in this result,

$$\text{Mean value of } dp_x dp_x = (p_{x^2} - p_x^2)/N \dots\dots\dots(19).$$

If we multiply (10) and (11), sum for all samples and divide by the number of samples we deduce

$$\begin{aligned} N\sigma_{xy}\sigma_{xz}R_{r_{xy},r_{xz}}/r_{xy}r_{xz} &= \frac{p_{xyz} - p_{xy}p_z}{p_{xy}p_z} - \frac{1}{2} \frac{(p_{xyz} - p_{xy}p_z)}{p_{xy}p_z} - \frac{1}{2} \frac{(p_{xyz} - p_{xy}p_z)}{p_{xy}p_z} \\ &- \frac{1}{2} \frac{(p_{x^2} - p_x^2)}{p_x^2} - \frac{1}{2} \frac{(p_{x^2} - p_x^2)}{p_x^2} + \frac{1}{4} \frac{(p_{x^2} - p_x^2)}{p_x^2} \\ &+ \frac{1}{4} \frac{(p_{x^2} - p_x^2)}{p_x^2} + \frac{1}{4} \frac{(p_{x^2} - p_x^2)}{p_x^2} + \frac{1}{4} \frac{(p_{x^2} - p_x^2)}{p_x^2} \dots\dots\dots(20). \end{aligned}$$

This result like Sheppard's formula for σ_{xy}^2 is much simpler when expressed in reduced moments. Let us write

$$\frac{p_x^2 y^m z^k}{\sigma_x^2 \sigma_y^m \sigma_z^k} = q_x^m y^m z^k,$$

so that q_x is unity and $q_{xy} = r_{xy}$. The numerical term in (20) is

$$-1 + \frac{1}{2}(4) + \frac{1}{4}(-4)$$

or zero, hence

$$\begin{aligned} N\sigma_{xy}\sigma_{xz}R_{r_{xy},r_{xz}}/r_{xy}r_{xz} &= \frac{q_{xyz}}{r_{xy}r_{xz}} - \frac{1}{2} \left(\frac{q_{xyz} + q_{xyz}}{r_{xy}} + \frac{q_{x^2z} + q_{x^2z}}{r_{xz}} \right) + \frac{1}{4} (q_{x^2z} + q_{x^2z} + q_{x^2z} + q_{x^2z}) \\ &\dots\dots\dots(21). \end{aligned}$$

In the same notation Sheppard's formula becomes

$$\frac{\sigma_{xy}^2}{r_{xy}^2} = \frac{1}{N} \left\{ \frac{q_{x^2y^2}}{r_{xy}^2} + \frac{1}{4} (\beta_3 + \beta'_3) + q_{x^2y^2} - \frac{q_{x^2y} + q_{xy^2}}{r_{xy}} \right\} \dots\dots\dots(22)*.$$

To find the correlation between r_{xy} and r_{xz} we have only to replace t by x in (21), thus

$$\begin{aligned} N\sigma_{xy}\sigma_{xz}R_{r_{xy},r_{xz}}/r_{xy}r_{xz} &= \frac{q_{x^2yz}}{r_{xy}r_{xz}} - \frac{1}{2} \left(\frac{q_{x^2yz} + q_{x^2yz}}{r_{xy}} + \frac{q_{x^2z} + q_{x^2z}}{r_{xz}} \right) + \frac{1}{4} (q_{x^2z} + q_{x^2z} + q_{x^2z} + q_{x^2z}) \\ &\dots\dots\dots(23). \end{aligned}$$

(3) These correlation coefficients will simplify if the regression be linear and simplify to a considerable extent if at the same time the distribution be normal.

For with linear regression

$$\begin{aligned} Np_{x^2yz} &= S S S (n_{xyz} x^2 y z) \dagger \\ &= S S (n_{xy} x^2 y \times \bar{x}_{xy}), \end{aligned}$$

where \bar{x}_{xy} is the mean value of z for given values of x and y .

* For the denominator of left-hand side, cf. *Biometrika*, Vol. IX. p. 4.

† The origin being taken at the mean.

But from the usual regression equation

$$\bar{z}_{xy} = x \left(\frac{r_{xz} - r_{xy}r_{yz}}{1 - r_{xy}^2} \right) \frac{\sigma_z}{\sigma_x} + y \left(\frac{r_{yz} - r_{xz}r_{xy}}{1 - r_{xy}^2} \right) \frac{\sigma_z}{\sigma_y},$$

so that for linear regression

$$p_{z^2xy} = \left(\frac{r_{xz} - r_{xy}r_{yz}}{1 - r_{xy}^2} \right) \frac{\sigma_z}{\sigma_x} p_{x^2y} + \left(\frac{r_{yz} - r_{xz}r_{xy}}{1 - r_{xy}^2} \right) \frac{\sigma_z}{\sigma_y} p_{y^2x} \dots\dots\dots(24).$$

Further

$$\begin{aligned} p_{x^2y} &= \frac{1}{N} S S (n_{xy} x^2 y) \\ &= \frac{1}{N} S_x (n_x x^2 \bar{y}_x) \\ &= \frac{1}{N} S_x \left(n_x x^2 \frac{\sigma_y}{\sigma_x} r_{xy} \right), \end{aligned}$$

or

$$p_{x^2y} = p_x r_{xy} \frac{\sigma_y}{\sigma_x} \dots\dots\dots(25),$$

while $q_{x^2y} = 1 + r_{xy}^2 \sqrt{(\beta_x - 1)(\beta_y - 1)}$ approximately(26)*,

so that (23) can be evaluated approximately by the use of simple moment coefficients and correlation coefficients only. *If in addition the distribution be normal, we know that*

$$p_{x^2y} = 3p_{xy}p_{x^2} \text{ and } p_{x^2y^2} = (1 + 2r_{xy}^2)/p_{xy}p_{y^2},$$

so that for normal distributions

$$\frac{p_{x^2yz}}{p_{xy}p_{yz}} = \left(\frac{r_{xz} - r_{xy}r_{yz}}{1 - r_{xy}^2} \right) \frac{3}{r_{xz}} + \frac{r_{yz} - r_{xz}r_{xy}}{1 - r_{xy}^2} \frac{1 + 2r_{xy}^2}{r_{xy}r_{xz}},$$

or

$$\frac{q_{x^2yz}}{r_{xy}r_{xz}} = 2 + r_{yz}/r_{xy}r_{xz} \dots\dots\dots(27).$$

Similarly

$$\frac{q_{xyz^2}}{r_{xy}} = 1 + \frac{2r_{yz}r_{xz}}{r_{xy}} \dots\dots\dots(28),$$

and

$$\frac{q_{xy^2z}}{r_{xz}} = 1 + 2r_{yz}r_{xy}/r_{xz} \dots\dots\dots(29),$$

Substituting these values in (23) we obtain, after some reduction and using

$$\begin{aligned} \sigma_{r_{xy}} &= (1 - r_{xy}^2)/\sqrt{N}, \\ (1 - r_{xy}^2)(1 - r_{xz}^2) R_{r_{xy}r_{xz}} &= r_{yz}(1 - r_{xy}^2)(1 - r_{xz}^2) - \frac{1}{2} r_{xy}r_{xz}(1 - r_{xy}^2 - r_{yz}^2 - r_{xz}^2 + 2r_{xy}r_{yz}r_{xz}) \\ &\dots\dots\dots(30), \end{aligned}$$

agreeing with (6) the value obtained by Pearson and Filon for normal distributions.

As regards the more general case dealing with the correlation between deviations of r_{xy} and those of r_{xz} given by equation (22), we have *when the regression is linear*

$$\begin{aligned} p_{xyzx} &= S S S S \{n_{xyzx} (xyzx)\} \\ &= S S S \{n_{xyz} (xyz \bar{r}_{xyz})\}, \end{aligned}$$

* *Biometrika*, Vol. IX. p. 4.

where \bar{t}_{xyz} is the mean value of t for given values of x, y and z , so that as is well known

$$\bar{t}_{xyz} = -x \frac{\sigma_t \Delta_{xt}}{\sigma_x \Delta_{tt}} - y \frac{\sigma_t \Delta_{yt}}{\sigma_y \Delta_{tt}} - z \frac{\sigma_t \Delta_{zt}}{\sigma_z \Delta_{tt}} \dots\dots\dots(31),$$

where

$$\Delta = \begin{vmatrix} 1, & r_{xy}, & r_{xz}, & r_{xt} \\ r_{xy}, & 1, & r_{yz}, & r_{yt} \\ r_{xz}, & r_{yz}, & 1, & r_{zt} \\ r_{xt}, & r_{yt}, & r_{zt}, & 1 \end{vmatrix} \dots\dots\dots(32),$$

and Δ_{xyt} is the minor corresponding to r_{xy} . Thus

$$q_{xyzt} = -(\Delta_{xt}q_{x^2z} + \Delta_{yt}q_{xy^2z} + \Delta_{zt}q_{xyz^2})/\Delta_{tt} \dots\dots\dots(33),$$

so that $R_{r_{xy}, r_{xz}}$ can be evaluated approximately in the case of linear regression without employing any mixed moments beyond the simple product moment occurring in a correlation coefficient.

For normal distributions we may use (27), (28) and (29) giving

$$q_{xyzt} = -[\Delta_{xt}(2r_{xy}r_{xz} + r_{yz}) + \Delta_{yt}(2r_{xy}r_{yz} + r_{xz}) + \Delta_{zt}(2r_{xz}r_{yz} + r_{xy})]/\Delta_{tt} \dots(34).$$

By well-known properties of first minors of a determinant we have from (32)

$$\Delta_{xt} + r_{xy}\Delta_{yt} + r_{xz}\Delta_{zt} + r_{xt}\Delta_{tt} = 0 \dots\dots\dots(35),$$

$$r_{xy}\Delta_{xt} + \Delta_{yt} + r_{yz}\Delta_{zt} + r_{yt}\Delta_{tt} = 0 \dots\dots\dots(36),$$

$$r_{xz}\Delta_{xt} + r_{yz}\Delta_{yt} + \Delta_{zt} + r_{zt}\Delta_{tt} = 0 \dots\dots\dots(37).$$

Multiply these equations by r_{yz}, r_{xz}, r_{xy} respectively and add,

$$\begin{aligned} \therefore (r_{yz} + 2r_{xz}r_{xy}) \Delta_{xt} + (r_{xz} + 2r_{yz}r_{xy}) \Delta_{yt} \\ + (r_{xy} + 2r_{yz}r_{xz}) \Delta_{zt} + (r_{yt}r_{xt} + r_{xz}r_{yt} + r_{xy}r_{zt}) \Delta_{tt} = 0 \dots\dots(38). \end{aligned}$$

Combining this result with (33) we see that for normal distributions

$$q_{xyzt} = r_{xy}r_{zt} + r_{yz}r_{xt} + r_{xz}r_{yt} \dots\dots\dots(39)*,$$

an interesting result likely to prove useful in other applications and probably capable of generalisation. Particular cases of (39) are obtained by putting $t = x$ so that $q_{x^2yz} = r_{yz} + 2r_{xy}r_{xz}$ which is (27) and $t = x, z = y$ giving $q_{x^2y^2} = 1 + 2r_{xy}^2$ which is well known.

If we now substitute these values in equation (21) we find

$$\begin{aligned} N\sigma_{r_{xy}}\sigma_{r_{xz}}R_{r_{xy}, r_{xz}} \\ = 2r_{yz}r_{xt} + 2r_{xz}r_{yt} - 2r_{xz}r_{yz}r_{xt} - 2r_{xt}r_{yt}r_{zt} \\ - 2r_{xz}r_{xt}r_{xy} - 2r_{yz}r_{yt}r_{xy} + r_{xy}r_{zt}(r_{xz}^2 + r_{xt}^2 + r_{yz}^2 + r_{yt}^2). \end{aligned}$$

The right-hand member can be put in the form

$$\begin{aligned} \frac{1}{2}\{(r_{xt} - r_{xz}r_{zt})(r_{yz} - r_{xy}r_{xz}) + (r_{xt} - r_{xy}r_{yt})(r_{yz} - r_{yt}r_{zt}) \\ + (r_{xz} - r_{xy}r_{yz})(r_{yt} - r_{yz}r_{zt}) + (r_{xz} - r_{xt}r_{yz})(r_{yt} - r_{xy}r_{zt})\}, \end{aligned}$$

* This result, which is accurate for normal distributions, is given as approximately true for such distributions by H. E. Soper, *Biometrika*, Vol. XL, p. 100

and if we remember that for normal distributions

$$\sigma_{r_{xy}} = (1 - r_{xy}^2)/\sqrt{N}, \quad \sigma_{r_{xz}} = (1 - r_{xz}^2)/\sqrt{N},$$

this result agrees with Pearson and Filon's value quoted above as equation (7).

(4) To find the probable error of a standard deviation.

$$\begin{aligned} \sigma_x^2 &= p_{xx}, \\ \therefore \frac{d\sigma_x}{\sigma_x} &= \frac{dp_{xx}}{2p_{xx}} \dots\dots\dots(40). \end{aligned}$$

Hence

$$\begin{aligned} \frac{\Sigma \sigma_x^2}{\sigma_x^2} &= \frac{p_{xx} - p_{xx}^2}{4p_{xx}^2 N} \text{ by (17)} \\ &= \frac{\beta_x - 1}{4N}, \\ \therefore \Sigma \sigma_x &= \frac{\sigma_x \sqrt{\beta_x - 1}}{2\sqrt{N}} \dots\dots\dots(41). \end{aligned}$$

This result is well known, and for normal distributions, i.e. when $\beta_x = 3$, becomes

$$\Sigma \sigma_x = \frac{\sigma_x}{\sqrt{2N}}, \text{ agreeing with (1).}$$

To find the correlation between a standard deviation σ_x and a correlation coefficient r_{yz} , we multiply (40) by the equation

$$\frac{dr_{yz}}{r_{yz}} = \frac{dp_{yz}}{p_{yz}} - \frac{dp_{yy}}{2p_{yy}} - \frac{dp_{zz}}{2p_{zz}},$$

and sum for all samples and divide by their number in the usual way, obtaining

$$\begin{aligned} 2N\sigma_x \sigma_{r_{yz}} R_{\sigma_x, r_{yz}} / \sigma_x r_{yz} & \\ &= (p_{xyyz} - p_{xx}p_{yy})/p_{xx}p_{yy} - \frac{1}{2}(p_{xyyz} - p_{xx}p_{yy})/p_{xx}p_{yy} - \frac{1}{2}(p_{xyyz} - p_{xx}p_{yy})/p_{xx}p_{yy} \\ &= \frac{q_{xyyz}}{r_{yz}} - \frac{q_{xyyz} + q_{xyyz}}{2} \dots\dots\dots(42), \end{aligned}$$

a result which as before can be approximated to in the case of linear regression, and which for normal distributions becomes*

$$\begin{aligned} R_{\sigma_x, r_{yz}} &= \frac{r_{yz} + 2r_{xy}r_{xz} - \frac{1}{2}(2 + 2r_{xy}^2 + 2r_{xz}^2)}{2N \cdot \frac{\sigma_x}{\sqrt{2N}} \cdot \frac{(1 - r_{yz}^2)}{\sqrt{N}} \cdot \frac{1}{\sigma_x r_{yz}}} \\ &= \frac{2r_{xy}r_{xz} - (r_{xy}^2 + r_{xz}^2)r_{yz}}{\sqrt{2} \cdot (1 - r_{yz}^2)} \dots\dots\dots(43), \end{aligned}$$

agreeing with equation (5).

* For the case $z=x$, i.e. $R_{\sigma_x, r_{xy}}$, cf. *Biometrika*, Vol. IX, p. 8.