" On the Significance of Bravais' Formulæ for Regression, &c.,
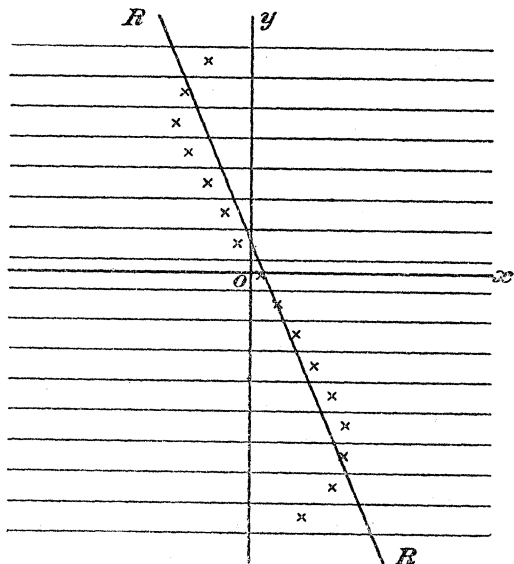    in the case of Skew Correlation." By G. UDNY YULE.
    Communicated by Professor KARL PEARSON, F.R.S.
    Received December 14, 1896,—Read February 18, 1897.

The only theory of correlation at present available for practical
use is based on the normal law of frequency, but, unfortunately, this
law is not valid in a great many cases which are both common and
important. It does not hold good, to take examples from biology,
for statistics of fertility in man, for measurements on flowers, or for
weight measurements even on adults. In economic statistics, on the
other hand, normal distributions appear to be highly exceptional:
variation of wages, prices, valuations, pauperism, and so forth, are
always skew. In cases like these we have at present no means of
measuring the correlation by one or more " correlation coefficients "
such as are afforded by the normal theory.

It seems worth while noting, under these circumstances, that in
ordinary practice statisticians never concern themselves with the
form of the correlation, normal or otherwise, but yet obtain results of
interest—though always lacking in numerical exactness and fre-
quently in certainty. Suppose the case to be one in which two
variables are varying together in time, curves are drawn exhibiting
the history of the two. If these two curves appear, generally
speaking, to rise and fall together, the variables are held to be corre-
lated. If on the other hand it is not a case of variation with time,
the associated pairs may be tabulated in order according to the
magnitude of one variable, and then it may be seen whether the
entries of the other variable also occur in order. Both methods are
of course very rough, and will only indicate very close correlation,
but they contain, it seems to me, the point of prime importance at
all events with regard to economic statistics. In all the classical
examples of statistical correlation (*e.g.*, marriage-rate and imports,
corn prices and vagrancy, out-relief and wages) we are only
primarily concerned with the question is a large $x$ usually associated
with a large $y$ (or small $y$); the further question as to the form of
this association and the relative frequency of different pairs of the
variables is, at any rate on a first investigation, of comparatively
secondary importance.

Let $Ox$, $Oy$ be the axes of a *three dimensional* frequency-surface
drawn through the mean $O$ of the surface parallel to the axes of
measurement, and let the points marked ⊗ be the means of succes-
sive $x$-arrays, lying on some curve that may be called the curve of
regression of $x$ on $y$. Now let a line, RR, be fitted to this curve,

subjecting the distances of the means from the line to some minimal condition. If the slope of RR is positive we may say that large values of $x$ are on the whole associated with large values of $y$, if it is negative large values of $x$ are associated with small values of $y$. Further, if the slope of RR to the vertical be given we shall have a measure of a rough practical kind of the shift of the mean of an $x$-array when its type $y$ is altered. The equation to RR consequently gives a concise and definite answer to two most important statistical questions. It is also evident that if the means of the arrays actually lie in a straight line (as in normal correlation), the equation to RR must be the equation to the line of regression.

Let $n$ be the number of observations in any $x$-array, and let $d$ be the horizontal distance of the mean of this array from the line RR. I propose to subject the line to the condition that the sum of all quantities like $nd^2$ shall be a minimum, *i.e.*, I shall use the condition of least squares. I do this solely for convenience of analysis; I do not claim for the method adopted any peculiar advantage as regards the probability of its results. It would, in fact, be absurd to do so, for I am postulating at the very outset that the curve of regression is only exceptionally a straight line; there can consequently be no meaning in seeking for the *most probable* straight line to represent the regression.

Let $x$, $y$ be a pair of associated deviations, let $\sigma$ be the standard deviation of any array about its own mean, and let

$$X = a + bY$$

be the equation to RR. Then for any one array

$$S\{x-(a+by)\}^2 = S\{x-(a+bY)\}^2 = n\sigma^2 + nd^2.$$

Hence, extending the meaning of S to summation over the whole surface

$$S(nd^2) = S\{x-(a+by)\}^2 - Sn\sigma^2.$$

But in this expression $S(n\sigma^2)$ is independent of $a$ and $b$, it is, in fact, a characteristic of the surface. Therefore, making $S(nd^2)$ a minimum is equivalent to making

$$S\{x-(a+by)\}^2$$

a minimum. That is to say, we may regard our method in another light. We may say that we form *a single-valued* relation

$$x = a + by$$

between a pair of associated deviations, such that the sum of the squares of our errors in estimating any one $x$ from its $y$ by the relation is a minimum. This single-valued relation, which we may call the characteristic relation, is simply the equation to the line of regression RR. There will be two such equations to be formed corresponding to the two lines of regression.

The idea of the method may at once be extended to the case of correlation between several variables $x_1$, $x_2$, $x_3$, &c. Let $n$ be the number of observations in an array of $x_1$'s associated with fixed values $X_2$, $X_3$, $X_4$, &c., of the remaining variables, let $\sigma_1$ be the standard deviation of this array, and let $d$ be the difference of its mean from the value given by a regression equation

$$X_1 = a_{12}X_2 + a_{13}X_3 + a_{14}X_4 + \ldots\ldots$$

Then, as before, we shall determine the coefficients $a_{12}$, $a_{13}$ $a_{14}$, &c., so as to make $Snd^2$ a minimum. But this is again equivalent to making

$$S\{x_1-(a_{12}x_2 + a_{13}x_3 + a_{14}x_4 + \ldots)\}^2$$

a minimum for

$$S\{x_1-(a_{12}x_2 + a_{13}x_3 + a_{14}x_4 + \ldots)\}^2 = S(n\sigma_1^2) + S(nd^2).$$

Hence, we may say that we solve for a single-valued relation

$$x_1 = a_{12}x_2 + a_{13}x_3 + a_{14}x_4 + \ldots$$

between our variables; the relation being such that the sum of the squares of the errors made in estimating $x_1$ from its associated values $x_2$, $x_3$, &c., is the least possible. In the case of normal correla-

tion this "characteristic relation" must become the "equation of regression" which gives the means of any $x_1$-array, as only in this way can $S nd^2$ be made a minimum, *i.e.*, zero.

It might be said that it would be more natural to form a "characteristic relation" between the absolute values of the variables and not their deviations from the mean. This may, however, be most conveniently done by working with the *mean* as origin until the characteristic is obtained, and then transferring the equation to zero as origin. It would be much more laborious and would only lead to the same result if zero were used *ab initio* as origin.

We may now proceed to the discussion of the special cases of two, three, or more variables. The actual formulæ obtained are not, it will be found, novel in themselves, but throw an unexpected light on the meaning of the expressions previously given by Bravais* for the case of normal correlation.

(1) *Case of Two Variables.*—Since $x$ and $y$ represent deviations from their respective means, we have, using S to denote summation over the whole surface,

$$S(x) = S(y) = 0.$$

The characteristic or regression equations which we have to find are of the form

$$\left. \begin{array}{l} x = a_1 + b_1 y \\ y = a_2 + b_2 x \end{array} \right\} \quad \dots\dots\dots\dots\dots\dots \quad (1).$$

Taking the equation for $x$ first, the normal equations for $a_1$ and $b_1$ are

$$\left. \begin{array}{l} S(x) = N a_1 + b_1 S(y) \\ S(xy) = a_1 S(y) + b_1 S(y^2) \end{array} \right\} \quad \dots\dots\dots\dots\dots \quad (2),$$

N being the total number of correlated pairs. From the first of these equations we have at once

$$a_1 = 0.$$

From the second

$$b_1 = \frac{S(xy)}{S(y^2)} .$$

To simplify our notation let us write

$$S(x^2) = N \sigma_1^2. \qquad Sy^2 = N \sigma_2^2.$$
$$S(xy) = N r \sigma_1 \sigma_2.$$

$\sigma_1$ and $\sigma_2$ are then the two standard-deviations or errors of mean

* "Mémoires par divers Savants," 1846, p. 255, and Professor Pearson's paper on "Regression, Heredity, &c." 'Phil. Trans.,' A, vol. 187 (1896), p. 261 *et seq.*

square. $r$ is Bravais' value of the coefficient of correlation. Re-writing $b_1$ in terms of these symbols, we have

$$b_1 = r \frac{\sigma_1}{\sigma_2} \dots\dots\dots\dots\dots\dots\dots \quad (3).$$

Similarly, $\qquad a_2 = 0, \qquad\quad b_2 = r \frac{\sigma_2}{\sigma_1} \dots\dots\dots\dots\dots \quad (4).$

But the expressions on the right of (3) and (4) are the values obtained by Bravais *on the assumption of normal correlation* for the regression of $x$ on $y$, and the regression of $y$ on $x$. That is to say, the Bravais values for the regressions are simply those values of $b_1$ and $b_2$, which make

$$S(x - b_1 y)^2 \text{ and } S(x - b_2 y)^2$$

respectively minima, *whatever be the form of the correlation between the two variables.* Again, whatever the form of the correlation, if the regression be really linear, the equations to the lines of regression are those given above (as we pointed out in the introduction). This theorem admits of a very simple and direct geometrical proof.

Let $n$ be the number of correlated pairs in any one array taken parallel to the axis of $x$, and let $\theta$ be the angle that the line of regression makes with the axis of $y$. Then, for a single array,

$$S(xy) = yS(x) = ny^2 \tan \theta,$$

or extending the significance of S to summation over the whole surface,

$$S(xy) = N \tan \theta \sigma_2{}^2,$$

that is,

$$\tan \theta = r \frac{\sigma_1}{\sigma_2}.$$

*In any case, then, where the regression appears to be linear, Bravais' formulæ may be used at once without troubling to investigate the normality of the distribution. The exponential character of the surface appears to have nothing whatever to do with the result.*

To return, again, to the most general case, we see that both coefficients of regression must have the same sign, namely, the sign of $r$. Hence, either regression will serve to indicate whether there is correlation or no, for there is no reason, *à priori*, why the values of $b_1$ and $b_2$, as determined above, should be positive rather than negative. But, nevertheless, the regressions are not convenient measures of correlation, for, on comparing two similar cases, we may find, say,

$$b_1 > b'_1, \qquad\qquad b_2 < b'_2,$$

where $b_1 b_2$, $b'_1 b'_2$ are the regressions in the two cases. To which distribution are we, in such a case, to attribute the greater correlation ? Bravais' coefficient solves the difficulty, we may say, in one way, by taking the geometrical mean of the two regressions as the measure of correlation. It will still remain valid for non-normal correlation. But there are other and less arbitrary interpretations even in the general case.

Suppose that instead of measuring $x$ and $y$ in arbitrary units we measure each in terms of its own standard deviation. Then let us write

$$\frac{x}{\sigma_1} = \rho \frac{y}{\sigma_2} \quad\dots\dots\dots\dots\dots\dots\dots \quad (5),$$

and solve for $\rho$ by the method of least squares. We have omitted a constant on the right-hand side, since it would vanish as before. We have, at once,

$$\rho = \frac{S(xy)}{S(y^2)} \frac{\sigma_2}{\sigma_1} = r \quad\dots\dots\dots\dots\dots\dots \quad (6),$$

That is to say, if we measure $x$ and $y$ each in terms of its own standard deviation, $r$ becomes at once the regression of $x$ on $y$, and the regression of $y$ on $x$. The regressions being, in fact, the fundamental physical quantities, $r$ is a coefficient of correlation because it is a coefficient of regression.[*]

Again, let us form the sums of the squares of residuals in equations (1) and (5). Inserting the values of $b_1$, $b_2$, and $\rho$, we have—

$$\left.\begin{aligned}
S(x - b_1 y)^2 &= N\sigma_1^2(1 - r^2) \\
S(y - b_2 x)^2 &= N\sigma_2^2(1 - r^2) \\
S\left(\frac{x}{\sigma_1} - \rho\frac{y}{\sigma_2}\right)^2 &= S\left(\frac{y}{\sigma_2} - \rho\frac{x}{\sigma_1}\right)^2 = N(1 - r^2)
\end{aligned}\right\} \dots\dots\dots \quad (7).$$

Any one of these quantities, being the sum of a series of squares, must be positive. Hence $r$ cannot be greater than unity. If $r$ be equal to unity, or if the correlation be perfect, all the above three sums become zero. But

$$S\left(\frac{x}{\sigma_1} + \frac{y}{\sigma_2}\right)^2$$

can only vanish if

$$\frac{x}{\sigma_1} + \frac{y}{\sigma_2} = 0$$

in every case, or if the relation hold good,

---

[*] That the regression becomes the coefficient of correlation when each deviation is measured in terms of its standard-deviation in the case of normal correlation has been pointed out by Mr. Francis Galton. *Vide* Pearson ' Phil. Trans.,' A, vol. 187, p. 307, note.

$$\frac{x_1}{y_1} = \frac{x_2}{y_2} = \frac{x_3}{y_3} = \ \dots \ = \pm \frac{\sigma_1}{\sigma_2} \ \dots\dots\dots\dots \quad (8),$$

the sign of the last term depending on the sign of $r$. Hence the statement that two variables are " perfectly correlated " implies that relation (8) holds good, or that all pairs of deviations bear the same ratio to one another. It follows that in correlation, where the means of arrays are not collinear, or the deviation of the mean of the array is not a linear function of the deviation of the type, $r$ can never be unity, though we know from experience that it can approach pretty closely to that value. If the regression be very far from linear, some caution must evidently be used in employing $r$ to compare two different distributions.

In the case of normal correlation, $\sigma_1 \sqrt{1-r^2}$ is the standard deviation of any array of the $x$ variables, corresponding to a single type of $y$'s. $\sigma_2 \sqrt{1-r^2}$ is similarly the standard deviation of any array of the $y$ variables, corresponding to a single type of $x$'s. In the general case, the first expression may be interpreted as the mean standard deviation of the $x$-arrays from the line of regression, and the second expression as the mean standard deviation of the $y$-arrays from the line of regression. Otherwise we may regard

$$\sigma_1 \sqrt{1-r^2}$$

as the standard error made in estimating $x$ from the relation

$$x = b_1 y,$$

and

$$\sigma_2 \sqrt{1-r^2}$$

as the standard error made in estimating $y$ from the relation

$$y = b_2 x,$$

these interpretations being independent of the form of the correlation.

### (2.) *Case of Three Variables.*

Let the three correlated variables be $X_1$, $X_2$, $X_3$, and let $x_1$, $x_2$, $x_3$ denote deviations of these variables from their respective means. Let us write, for brevity,

$$S(x_1^2) = N\sigma_1^2, \qquad S(x_2^2) = N\sigma_2^2$$

$$S(x_3^2) = N\sigma_3^2$$

$$S(x_1 x_2) = N r_{12} \sigma_1 \sigma_2$$

$$S(x_2 x_3) = N r_{23} \sigma_2 \sigma_3$$

$$S(x_3 x_1) = N r_{31} \sigma_3 \sigma_1.$$

Our characteristic or regression-equation will now be of the form

$$x_1 = b_{12}x_2 + b_{13}x_3 \dots \dots \dots \dots \dots \quad (9),$$

$b_{12}$ and $b_{13}$ being the unknowns to be determined from the observations by the method of least squares. I have omitted a constant term on the right-hand side, since its least-square value would be zero as before. The two normal equations are now—

$$S(x_1x_2) = b_{12}S(x_2^2) + b_{13}S(x_2x_3)$$

$$S(x_1x_3) = b_{12}S(x_2x_3) + b_{13}S(x_3^2),$$

or replacing the sums by the symbols defined above, and simplifying—

$$\left. \begin{array}{l} r_{12}\sigma_1 = b_{12}\sigma_2 + b_{13}r_{23}\sigma_3 \\ r_{13}\sigma_1 = b_{12}r_{23}\sigma_2 + b_{13}\sigma_3 \end{array} \right\} \dots \dots \dots \dots \quad (10),$$

whence

$$\left. \begin{array}{l} b_{12} = \dfrac{r_{12} - r_{13}r_{23}}{1 - r_{23}^2} \ \dfrac{\sigma_1}{\sigma_2} \\[2mm] b_{13} = \dfrac{r_{13} - r_{12}r_{23}}{1 - r_{23}^2} \ \dfrac{\sigma_1}{\sigma_3} \end{array} \right\} \dots \dots \dots \dots \quad (11).$$

That is, the characteristic relation between $x_1$ and $x_2 x_3$ is—

$$x_1 = \frac{r_{12} - r_{13}r_{23}}{1 - r_{23}^2} \frac{\sigma_1}{\sigma_2} x_2 + \frac{r_{13} - r_{12}r_{23}}{1 - r_{23}^2} \frac{\sigma_1}{\sigma_3} x_3 \dots \dots \dots \quad (12).$$

Now Bravais showed that *if the correlation were normal*, and we selected a group or array of $X_1$'s with regard to special values $h_2$ and $h_3$ of $x_2$ and $x_3$, then $h_1$ being the deviation of the mean of the selected $X_1$'s from the $X_1$-mean of the whole material,

$$h_1 = b_{12}h_2 + b_{13}h_3,$$

where $b_{12}$ and $b_{13}$ have the values given in (11). But evidently the relation is of much greater generality; it holds good so long as $h_1$ is a linear function of $h_2$ and $h_3$, *whatever be the law of frequency.*

Further, the values of $b_{12}$ and $b_{13}$ above determined, are, under any circumstances, such that

$$Sv^2 = S[x_1 - (b_{12}x_2 + b_{13}x_3)]^2,$$

is a minimum. If we insert in this expression the values of $b_{12}$ and $b_{13}$ from (11), we have, after some reduction,

$$S(v^2) = N\sigma_1^2 \left\{ 1 - \frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{23}r_{31}}{1 - r_{23}^2} \right\}$$

$$= N\sigma_1^2 \{ 1 - R_1^2 \} \dots \dots \dots \dots \dots \dots \quad (13),$$

say. In normal correlation $\sigma_1 \sqrt{1 - R_1^2}$ is the standard deviation of an $X_1$-array, corresponding to any given types of $X_2$ and $X_3$. In general correlation it may be regarded as the mean standard deviation of the $X_1$-arrays from the plane

$$x_1 = b_{12}x_2 + b_{13}x_3,$$

or as the standard error made in estimating $x_1$ from $x_2$ and $x_3$ by relation (12).

The quantity R is of some interest, as it exactly takes the place of $r$ in the residual expressions (7). $R_1$ may, in fact, be regarded as a coefficient of correlation between $x_1$ and $(x_2 x_3)$; it can only be unity if the linear relation (9) or (12) hold good in every case.

The quantities $b_{12}$, $b_{13}$, &c. (the others may be written down by symmetry), may be termed the net regressions of $x_1$ on $x_2$, $x_1$ on $x_3$, &c. If we write 2 for 1 and 1 for 2 in the value of $b_{12}$, we have

$$b_{21} = \frac{r_{12} - r_{13}r_{23}}{1 - r_{23}^2} \frac{\sigma_2}{\sigma_1},$$

$b_{21}$ being the the net regression of $x_2$ on $x_1$. In normal correlation, $b_{12}$ and $b_{21}$ are the regressions for any group of $X_1$'s or $X_2$'s associated with a fixed type of $X_3$'s. Hence, in this case (normal correlation), the coefficient of correlation for such a group is the geometrical mean of the two regressions, or

$$\rho_{12} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}}$$

a quantity that may be called the net coefficient of correlation between $x_1$ and $x_2$.* The similar net coefficients between $x_1$ and $x_3$, $x_2$ and $x_3$, may be written down by interchanging the suffixes.

In normal correlation $\rho_{12}$ is quite strictly the coefficient of correlation for *any* sub-group of $X_1$'s and $X_2$'s, whatever the associated type of $X_3$'s. In generalised correlation this will not be so, and $\rho_{12}$ can only retain an average significance.

The method does not appear to be capable of investigating changes in the net coefficient as we pass from one type to another, but it may be noted that whatever the form of the correlation, $\rho_{12}$ retains three of the chief properties of the ordinary coefficients : (1) it can only be

---

* My quantities, $b_{12}$, $b_{13}$, &c., were termed by Professor Pearson (" Regression &c.," ' Phil. Trans.,' A, vol. 187 (1896), p. 287), "Coefficients of double regression," and quantities like $b_{12} \frac{\sigma_2}{\sigma_1}$, $b_{13} \frac{\sigma_3}{\sigma_1}$, &c., "coefficients of double correlation." My quantities $\rho$ he did not use. Having named the $\rho$'s "net correlation," it seemed most natural to rename the $b$'s " net regressions," as the $b$'s and $\rho$'s are corresponding quantities.

Some of my results given above were quoted by Professor Pearson in his paper (*loc. cit.*, notes on pp. 268 and 287).

zero if both net regressions are zero; (2) it is a symmetrical function of the variables; (3) it cannot be greater than unity; for, by (13),

$$r_{12}^2 + r_{13}^2 - 2r_{12}r_{23}r_{31} < 1 - r_{23}^2,$$

or adding $r_{13}^2 r_{23}^2$ to both sides, and transferring $r_{13}^2$ to the right-hand side

$$(r_{12} - r_{13}r_{23})^2 < (1 - r_{13}^2)(1 - r_{23}^2).$$

If any two coefficients, say $r_{12}r_{13}$, be supposed known, the inequality we have used above will give us limits for the value of the third. Throwing it into the form

$$(r_{23} - r_{12}r_{13})^2 < 1 + r_{12}^2 r_{13}^2 - r_{12}^2 - r_{13}^2,$$

we have $r_{23}$ must lie between the limits

$$r_{12}r_{13} \pm \sqrt{r_{12}^2 r_{13}^2 - r_{12}^2 - r_{13}^2 + 1}.$$

The values of these limits for some special cases are collected in the following table :—

| Values of $r_{12}$ and $r_{13}$. | Limits of $r_{23}$. |
|---|---|
| $r_{12} = r_{13} = 0$ | 0 |
| $r_{12} = r_{13} = \pm 1$ | $+1$ |
| $r_{12} = +1,\ r_{13} = -1$ | $-1$ |
| $r_{12} = 0,\ r_{13} = \pm 1$ | 0 |
| $r_{12} = 0,\ r_{13} = \pm r$ | $\pm \sqrt{1 - r^2}$ |
| $r_{12} = r_{13} = \pm r$ | 1 and $2r^2 - 1$ |
| $r_{12} = +r,\ r_{13} = -r$ | $2r^2 - 1$ and $-1$ |
| $r_{12} = r_{13} = \pm \sqrt{0.5} = 0.707$ | 0 and 1 |
| $r_{12} = + \sqrt{0.5}\ \ r_{12} = -\sqrt{0.5}$ | 0 ,, $-1$ |

One is rather prone to argue that if A be correlated with B, and B with C, A will be correlated with C. Evidently this is not necessary. A may be positively correlated with B, and B positively correlated with C, but yet A may, in general, be negatively correlated with C. Only, if the coefficients (AB) and (BC) are both numerically greater than 0·707, can one even ascribe the correct sign to the (AC) correlation.

It is evident that one would, in general, expect to make a smaller standard error in estimating $x_1$ from the two associated variables $x_2$ and $x_3$, than in estimating it from one only, say $x_2$. But it seems desirable to prove this specifically, and to investigate under what conditions it will hold good. The necessary condition is—

$$\frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{23}r_{13}}{1 - r_{23}^2} > r_{12}^2,$$

that is,

$$r_{12}^2 + r_{13}^2 - 2r_{12}r_{13}r_{23} > r_{12}^2 - r_{12}^2 r_{13}^2,$$

or

$$(r_{13} - r_{12}r_{23})^2 > 0.$$

But $(r_{13} - r_{12}r_{23})$ is the numerator of $\rho_{13}$, the net coefficient of correlation between $x_1$ and $x_3$. Hence the standard error in the second case will be always less than in the first, so long as $\rho_{13}$ is not zero. The condition is somewhat interesting.

To take an arithmetical example, suppose one had in some actual case

$$r_{12} = +0\cdot8$$
$$r_{23} = +0\cdot5 \qquad r_{13} = +0\cdot4.$$

One might very naturally imagine that the introduction of the third variable with a fairly high correlation coefficient ($0\cdot4$) would considerably lessen the standard deviation of the $x_1$-array; but this is not so, for

$$\rho_{13} = \frac{0\cdot4 - (0\cdot5 \times 0\cdot8)}{\sqrt{0\cdot75 \times 0\cdot36}} = 0,$$

so the third variable would be of no assistance.

### III. *Case of Four Variables.*

This case is, perhaps, of sufficient practical importance to warrant our developing the results at length as in the last.

If $x_1$, $x_2$, $x_3$, $x_4$, be the associated deviations of the four variables from their respective means, the characteristic equation will be of the form

$$x_1 = b_{12}x_2 + b_{13}x_3 + b_{14}x_4 \dots\dots\dots\dots (14).$$

The normal equations for the $b$'s are, in our previous notation,

$$\left.\begin{array}{l} r_{12}\sigma_1 = b_{12}\sigma_2 + b_{13}r_{23}\sigma_3 + b_{14}r_{24}\sigma_4 \\ r_{13}\sigma_1 = b_{12}r_{23}\sigma_2 + b_{13}\sigma_3 + b_{14}r_{34}\sigma_4 \\ r_{14}\sigma_1 = b_{12}r_{24}\sigma_2 + b_{13}r_{34}\sigma_3 + b_{14}\sigma_4 \end{array}\right\}$$

Hence

$$b_{12} = \frac{\begin{vmatrix} r_{12} & r_{23} & r_{24} \\ r_{13} & 1 & r_{31} \\ r_{14} & r_{34} & 1 \end{vmatrix}}{\begin{vmatrix} 1 & r_{23} & r_{24} \\ r_{23} & 1 & r_{34} \\ r_{24} & r_{34} & 1 \end{vmatrix}} \cdot \frac{\sigma_1}{\sigma_2} \dots\dots\dots\dots (15),$$

and so on for the others, $b_{12}$, $b_{13}$, &c., we may call the net regressions of $x_1$ on $x_2$, $x_1$ on $x_3$, &c., as before. By parity of notation, we have

$$b_{21} = \frac{\begin{vmatrix} r_{12} & r_{23} & r_{24} \\ r_{13} & 1 & r_{34} \\ r_{14} & r_{34} & 1 \end{vmatrix}}{\begin{vmatrix} 1 & r_{13} & r_{14} \\ r_{13} & 1 & r_{34} \\ r_{14} & r_{34} & 1 \end{vmatrix}} \frac{\sigma_2}{\sigma_1}$$

and we may again call

$$\rho_{12} = \sqrt{b_{12}b_{21}},$$

the net coefficient of correlation between $x_1$ and $x_2$. Expanding the determinants, we have, in fact,

$$\rho_{12} = \frac{r_{12}(1-r_{34}^2) + r_{13}(r_{34}r_{24}-r_{23}) + r_{14}(r_{23}r_{34}-r_{24})}{\sqrt{[(1-r_{34}^2) + r_{23}(r_{34}r_{24}-r_{23}) + r_{24}(r_{23}r_{34}-r_{24})][(1-r_{34}^2) + r_{13}(r_{34}r_{14}-r_{13}) + r_{14}(r_{13}r_{34}-r_{14})]}}$$

$$\dots\dots (16).$$

There are six such net coefficients, $\rho_{12}$, $\rho_{13}$, $\rho_{14}$, $\rho_{23}$, $\rho_{24}$, $\rho_{34}$. The above values of the regressions are again those usually obtained on the assumption of normal correlation.* The net correlation $\rho_{12}$ becomes, on that assumption, the coefficient of correlation for any group of the $x_1$ $x_2$ variables associated with fixed types of $x_3$ and $x_4$. If we write

$$u = x_1 - (b_{12}x_2 + b_{13}x_3 + b_{14}x_4),$$

we have, after some rather lengthy reduction,

$$\frac{1}{N} S(u^2) = \sigma_1^2(1 - R_1^2),$$

where

$$R_1^2 = \frac{\left\{ \begin{array}{c} r_{12}^2 + r_{13}^2 + r_{14}^2 - r_{12}^2 r_{34}^2 - r_{23}^2 r_{14}^2 - r_{13}^2 r_{24}^2 \\ -2(r_{13}r_{14}r_{34} + r_{12}r_{14}r_{24} + r_{12}r_{13}r_{23}) + 2(r_{12}r_{14}r_{23}r_{34} + r_{13}r_{14}r_{23}r_{24} + r_{12}r_{13}r_{24}r_{34}) \end{array} \right\}}{1 - r_{23}^2 - r_{34}^2 - r_{24}^2 + 2r_{23}r_{34}r_{24}}.$$

In normal correlation, $\sigma_1\sqrt{1 - R_1^2}$ is the standard deviation of all $x_1$-arrays associated with fixed types of $x_2$, $x_3$, and $x_4$. In general correlation, it is most easily interpreted as the standard error made in estimating $x_1$, by equation (14), from its associated values of $x_2$, $x_3$, and $x_4$.

As in the case of three variables, the quantity R may be considered as a coefficient of correlation. It can range between $\pm 1$, and can only become unity if the linear relation (14) hold good in each individual instance.

We showed at the end of the last section that the standard error made in estimating $x_1$ from the relation

$$x_1 = b_{12}x_2 + b_{13}x_3$$

* Professor Pearson, " Regression, Heredity, and Panmixia." ' Phil. Trans.,'
A, vol. 187 (1896), p. 294.

was always less than the standard error when only $x_2$ was taken into account, unless

$$\rho_{13} = 0.$$

We may now prove the similar theorem that when we use three variables, $x_2$, $x_3$, $x_4$, on which to base the estimate, the standard error will be again decreased, unless

$$\rho_{14} = 0.$$

The condition that $S(u^2)$, in our present case, shall be less than $S(r^2)$ in the last, is, in fact,

$$\left\{ \begin{array}{l} r_{12}{}^2 + r_{13}{}^2 + r_{14}{}^2 - r_{12}{}^2 r_{34}{}^2 - r_{23}{}^2 r_{14}{}^2 - r_{13}{}^2 r_{24}{}^2 \\ -2(r_{13}r_{14}r_{34} + r_{12}r_{13}r_{23} + r_{12}r_{14}r_{24}) \\ +2(r_{12}r_{14}r_{23}r_{34} + r_{14}r_{13}r_{24}r_{23} + r_{12}r_{13}r_{24}r_{34}) \end{array} \right\} (1 - r_{23}{}^2)$$
$$> (r_{12}{}^2 + r_{13}{}^2 - 2r_{12}r_{13}r_{23})(1 - r_{23}{}^2 - r_{24}{}^2 - r_{34}{}^2 + 2r_{23}r_{24}r_{34}).$$

This may be finally reduced to—

$$(r_{14} - r_{13}r_{34} - r_{12}r_{24} - r_{14}r_{23}{}^2 + r_{13}r_{23}r_{24} + r_{12}r_{23}r_{34})^2 > 0,$$

that is $$\rho_{14}{}^2 > 0.$$

The treatment of the general case of $n$ variables, so far as regards obtaining the regressions, is obvious, and it is unnecessary to give it at length.

We can now see that the use of normal regression formulæ is quite legitimate in all cases, so long as the necessary limitations of interpretation are recognised. Bravais' $r$ always remains a coefficient of correlation. These results I must plead as justification for my use of normal formulæ in two cases* where the correlation was markedly non-normal.

" Mathematical Contributions to the Theory of Evolution.—On a Form of Spurious Correlation which may arise when Indices are used in the Measurement of Organs." By KARL PEARSON, F.R.S., University College, London. Received December 29, 1896,—Read February 18, 1897.

(1) If the ratio of two absolute measurements on the same or different organs be taken it is convenient to term this ratio an *index*.

If $u = f_1(x, y)$ and $v = f_2(z, y)$ be two functions of the three variables $x$, $y$, $z$, and these variables be selected at random so that there exists no correlation between $x,y$, $y,z$, or $z,x$, there will still be found to

---

* 'Economic Journal,' Dec., 1895, and Dec., 1896, " On the Correlation of Total Pauperism with Proportion of Out-relief."