

## ON THE FREQUENCY DISTRIBUTION OF CERTAIN RATIOS

BY H. L. RIETZ

University of Iowa

Considerable interest in the distribution of ratios,  $t = y/x$ , has no doubt been suggested by important applications. For example, we may mention the opsonic index in bacteriology, the ratio of systolic to diastolic blood pressure in physiology, and ratios such as link relatives or certain index numbers in economics.

In 1910, Karl Pearson<sup>1</sup> gave certain properties of the distribution of ratios by means of approximate formulas for moments up to order four in terms of means, variances, product moments, and coefficients of variability of  $x$  and  $y$ . The resulting formulas did not give, with sufficient accuracy, the constants of the distribution of the opsonic index for the purpose of Dr. Greenwood to whom Pearson attributed the derivation of the formulas for the special case in which  $x$  and  $y$  are uncorrelated. Pearson next adopted the plan of tabulating the reciprocals, say  $x' = \frac{1}{x}$ , and then finding the constants of the distribution of the product  $yx'$  in the case in which  $x'$  and  $y$  are uncorrelated. He then obtained satisfactory results in illustrative examples.

In 1929, C. C. Craig<sup>2</sup> obtained the semi-invariants of  $y/x$  in terms of moments of  $x$  and  $y$ , and then expressed the moments in terms of the semi-invariants of the distribution function,  $f(x, y)$ , of  $x$  and  $y$ . By this means, he was able to deal with the case in which  $x$  and  $y$  are normally correlated under suitable conditions. Craig found it desirable to restrict the distribution of  $x$  in such a way that the probability of a zero value of  $x$  is an infinitesimal of sufficiently high order that a certain integral exists. This limitation seems to imply in applications to actual data that no zero values of  $x$  are to occur. This suggests that we deal with the cases of  $x$  at or near zero with considerable care.

By starting with the assumption that the values of  $x$  and  $y$  are a set of normally distributed pairs of values with correlation coefficient  $r$ , and by considering the quotient  $z = \frac{b + y}{a + x}$ ,  $a$  and  $b$  being constants, R. C. Geary,<sup>3</sup> in a paper published in 1930, found an algebraic function,  $u = f(z)$ , of fairly simple form with the property that  $u$  is nearly normally distributed with arithmetic mean zero and standard deviation unity provided that  $a + x$  is unlikely to

<sup>1</sup> On the constants of index distributions, *Biometrika*, Vol. 7 (1910), pp. 531-546.

<sup>2</sup> The frequency function of  $y/x$ , *Annals of Mathematics*, Vol. 30 (1928-29), pp. 471-486.

<sup>3</sup> The frequency distribution of the quotient of two normal variates, *J. Royal Statistical Society*. Vol. XCIII (1930), pp. 442-7.

have negative values. Here we have again a suggestion to exercise special care in the case of quotients with the divisor near zero or negative.

In 1932, Fieller<sup>4</sup> obtained in explicit form the approximate distribution of  $t = y/x$  where values  $(x, y)$  are drawn from the bivariate normal distribution

$$\frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-r^2}} e^{-\frac{1}{2} \frac{1}{1-r^2} \left\{ \frac{(x-\bar{x})^2}{\sigma_x^2} + \frac{(y-\bar{y})^2}{\sigma_y^2} - 2r \frac{(x-\bar{x})(y-\bar{y})}{\sigma_x\sigma_y} \right\}}$$

under the condition that  $\bar{x}$  is large compared with  $\sigma_x$ .

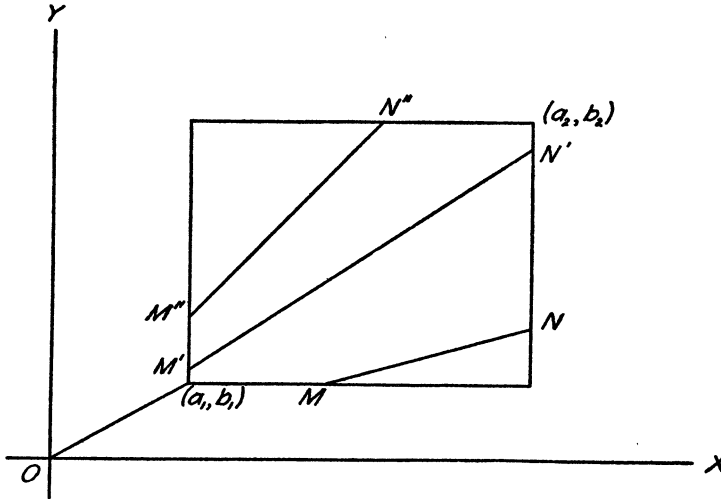


FIG. 1

Very recently Kullback<sup>5</sup> found the distribution law of the quotient,  $t = y/x$ , where  $x$  and  $y$  are drawn from Pearson Type III parent populations given by

$$f_1(x) = \frac{e^{-x} x^{p-1}}{\Gamma(p)}; \quad f_2(y) = \frac{e^{-y} y^{q-1}}{\Gamma(q)}, \quad 0 \leq x \leq \infty, \quad 0 \leq y \leq \infty.$$

It is fairly easy to see, in a general way, that the distribution of  $t = y/x$  depends very much on the location of the origin as well as on the parent distribution from which  $x$  and  $y$  are drawn. This fact will be fairly obvious from the present paper whose main purpose is to give clear geometrical descriptions of the distributions of ratios,  $t = y/x$ , for each of several cases in which  $(x, y)$  are points taken at random from certain simple geometrical figures conveniently located with respect to the origin.

In accord with the suggestions to be cautious when the divisor is near zero or negative, we consider first the very simple case of ratios  $t = y/x$  obtained

<sup>4</sup> E. C. Fieller, The distribution of the index in a normal bivariate population, *Biometrika*, Vol. 24 (1932), pp. 428-440.

<sup>5</sup> Solomon Kullback, *Annals of Mathematical Statistics*, Vol. VII (1936), pp. 51-53.

from points uniformly distributed over a rectangle such as is shown in Fig. 1 with sides parallel to coordinate axes and  $a_1 > 0, b_1 > 0$ . As indicated on Fig. 1, we assume for simplicity that the coordinates of the points are positive and  $a_1 \leq x \leq a_2, b_1 \leq y \leq b_2$ .

Case I. When  $\frac{b_1}{a_1} \leq \frac{b_2}{a_2}$ , Fig. 1.

Let  $k \, dx \, dy$  be the probability that a point  $(x, y)$  taken at random in the rectangle will fall into  $dxdy$  where  $k$  is a constant. Then

$$k \int_{b_1}^{b_2} \int_{a_1}^{a_2} dxdy = k(a_2 - a_1)(b_2 - b_1) = 1,$$

and 
$$k = \frac{1}{(a_2 - a_1)(b_2 - b_1)}.$$

Transform the element  $k \, dxdy$  into one with variables  $t$ , and  $x$  by making

$$x = x,$$

$$y = tx.$$

The Jacobian is  $|x| = x$ .

The new element is  $k \, x \, dxdt$  and is to be integrated over the range on  $x$  for an assigned  $t$  in order to get the probability, to within infinitesimals of higher order, that a random  $t$  falls into an assigned  $dt$ . By assigning  $t$  any value such that  $\frac{b_1}{a_2} \leq t \leq \frac{b_1}{a_1}$ , say  $t$  is the slope of  $MN$ , (Fig. 1), we have

$$(1) \quad k \int_{\frac{b_1}{a_2}}^{\frac{b_1}{a_1}} x dx dt = \frac{k}{2} \left( a_2^2 - \frac{b_1^2}{t^2} \right) dt$$

the limits of integration being indicated by the ends of the line  $MN$ .

When the assigned  $t$  is such that  $\frac{b_1}{a_1} \leq t \leq \frac{b_2}{a_2}$ , say  $t$  is the slope of the line  $M'N'$ , we have

$$(2) \quad k \int_{a_1}^{a_2} x dx dt = \frac{k}{2} (a_2^2 - a_1^2) dt$$

When the assigned  $t$  is such that  $\frac{b_2}{a_2} \leq t \leq \frac{b_2}{a_1}$ , say it is the slope of  $M''N''$ , we have

$$(3) \quad k \int_{a_1}^{\frac{b_2}{t}} x dx dt = \frac{k}{2} \left( \frac{b_2^2}{t^2} - a_1^2 \right) dt$$

Thus, from (1), (2), (3), when as in Fig. 1,  $\frac{b_1}{a_1} \leq \frac{b_2}{a_2}$ , the frequency function of  $t$  is given by

$$(4) \quad F(t) = \frac{k}{2} \left( a_2^2 - \frac{b_1^2}{t^2} \right) \quad \text{when} \quad \frac{b_1}{a_2} \leq t \leq \frac{b_1}{a_1},$$

$$(5) \quad F(t) = \frac{k}{2} (a_2^2 - a_1^2) \quad \text{when} \quad \frac{b_1}{a_1} \leq t \leq \frac{b_2}{a_2},$$

$$(6) \quad F(t) = \frac{k}{2} \left( \frac{b_2^2}{t^2} - a_1^2 \right) \quad \text{when} \quad \frac{b_2}{a_2} \leq t \leq \frac{b_2}{a_1}.$$

See Fig. 2 for the general form of the frequency curve  $F(t)$  when  $\frac{b_1}{a_1} < \frac{b_2}{a_2}$  with the segment from  $t = \frac{b_1}{a_1}$  to  $\frac{b_2}{a_2}$  a horizontal straight line and with discontinuities in the first derivatives of  $F(t)$  at  $t = \frac{b_1}{a_1}$  and  $t = \frac{b_2}{a_2}$ .

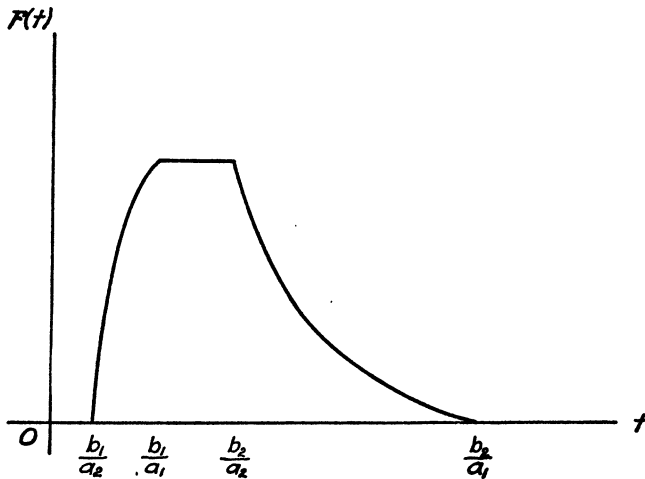


FIG. 2

When  $a_1 \rightarrow 0$ , and  $b_1 = 0$ , the frequency curve approaches

$$(7) \quad F(t) = \frac{a_2}{2b_2} \quad \text{when} \quad 0 \leq t \leq \frac{b_2}{a_2}$$

$$(8) \quad F(t) = \frac{b_2}{2a_2 t^2} \quad \text{when} \quad t \geq \frac{b_2}{a_2}.$$

It may be noted that the curve given by making  $a_1 = 0$  and  $b_1 = 0$  extends to infinity, and that the first and second moments about the origin are each infinite.

Case II. When  $\frac{b_1}{a_1} > \frac{b_2}{a_2}$ .

If the rectangle in Fig. 1 were moved upward keeping its sides parallel to the  $x$  and  $y$  axes until  $\frac{b_1}{a_1} > \frac{b_2}{a_2}$ , we would obtain

$$(9) \quad F(t) = \frac{k}{2} \left( a_2^2 - \frac{b_1^2}{t^2} \right) \quad \text{if} \quad \frac{b_1}{a_2} \leq t \leq \frac{b_2}{a_2},$$

$$(10) \quad F(t) = \frac{k}{2t^2} (b_2^2 - b_1^2) \quad \text{if} \quad \frac{b_2}{a_2} \leq t \leq \frac{b_1}{a_1},$$

$$(11) \quad F(t) = \frac{k}{2} \left( \frac{b_2^2}{t^2} - a_1^2 \right) \quad \text{if} \quad \frac{b_1}{a_1} \leq t \leq \frac{b_2}{a_1}.$$

By comparing (5) and (10), it may be observed that  $F(t)$  of the middle segment of the distribution curve differs much in Case II from its corresponding constant value in Case I.

By moving the rectangle of Fig. 1 downward, keeping its sides parallel to the  $x$  and  $y$  axes until  $b_1$  is negative, we easily find further forms of the distribution curve  $F(t)$ .

To consider the distribution of the ratio  $t = y/x$  for another very simple type of distribution of  $x$  and  $y$ , suppose we have given the distribution function

$$(12) \quad f(x, y) = k e^{-\frac{x}{a} - \frac{y}{b}}, \quad \left( \begin{array}{l} x \geq c > 0, y \text{ non-negative} \\ a > c, b > 0 \end{array} \right)$$

where  $\int_0^\infty \int_c^\infty f(x, y) dx dy = 1$ . Then

$$k = \frac{e^{c/a}}{ab}.$$

In this case,

$$(13) \quad \begin{aligned} F(t) &= \frac{e^{c/a}}{ab} \int_c^\infty x e^{-\frac{x}{a} - \frac{x t}{b}} dx \\ &= \frac{1}{b + at} \left( c + \frac{ab}{b + at} \right) e^{-\frac{c t}{b}}, \end{aligned}$$

a monotone decreasing function from  $t = 0$  to  $t = \infty$ .

With  $c = 0$  as a limiting value, we obtain

$$(14) \quad F(t) = \frac{ab}{(b + at)^2},$$

a distribution curve with the mean value of  $t$  at infinity.

If we should similarly consider

$$(15) \quad f(x, y) = \frac{2}{\pi\sigma_x\sigma_y} e^{-\frac{x^2}{2\sigma_x^2} - \frac{y^2}{2\sigma_y^2}} \quad (x \text{ and } y \text{ non-negative})$$

we easily obtain

$$(16) \quad F(t) = \frac{1}{2\pi\sigma_x\sigma_y \left( \frac{1}{\sigma_x^2} + \frac{t^2}{\sigma_y^2} \right)}$$

as the distribution function.

Although the difficulties<sup>6</sup> of the problem of the distribution of the ratio  $y/x$  when  $x$  and  $y$  are normally correlated have been overcome<sup>7</sup> to a considerable

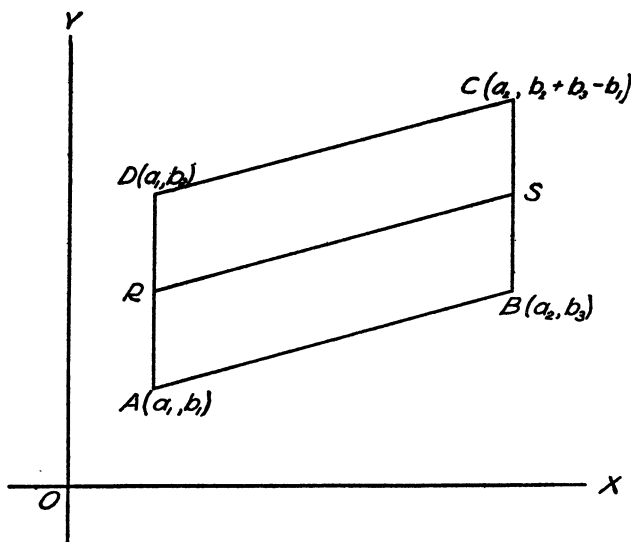


FIG. 3

extent, still the examination of some very simple geometric cases of non-normal but linear correlation may not be without some interest. Such a case will now be considered.

For one very simple case in which  $x$  and  $y$  are correlated, suppose we are given a set of points  $(x, y)$  uniformly distributed over the parallelogram  $ABCD$  (Fig. 3) with sides  $AD$  and  $BC$  parallel to the  $y$ -axis so that the regression of  $y$  on  $x$  is linear as shown by the line  $RS$ .

The equation of  $RS$  is

$$(17) \quad y = m(x - a_1) + \frac{b_1 + b_2}{2}.$$

<sup>6</sup> Loc. cit., Pearson, p. 531.

<sup>7</sup> Loc. cit., C. C. Craig, R. C. Geary, E. C. Fieller.

Then although  $x_i$  and  $y_i$  are correlated,  $x_i$  and

$$y'_i = y_i - m(x_i - a_i) - \frac{b_1 + b_2}{2}$$

are uncorrelated. Let us consider the distribution of the ratio  $t' = \frac{y'_i}{x_i}$ .

Consider the element of frequency  $k dx dy'$ , where

$$(18) \quad k(b_2 - b_1)(a_2 - a_1) = 1.$$

Change variables to  $x$  and  $t'$  by the transformation

$$\begin{aligned} x &= x, \\ y' &= t'x. \end{aligned}$$

Then the element of frequency becomes

$$(19) \quad kx dx dt'.$$

Next integrate (19) with respect to  $x$  under the restriction that  $t'$  is assigned. Three cases occur:

(a) When  $-\frac{b_2 - b_1}{2a_2} \leq t' \leq \frac{b_2 - b_1}{2a_2}$ , we obtain by integration of (19) for the element of relative frequency of  $t'$  in  $dt'$ ,

$$(20) \quad k \int_{a_1}^{a_2} x dx dt' = \frac{k}{2} (a_2^2 - a_1^2) dt'.$$

(b) When  $t' \geq \frac{b_2 - b_1}{2a_2}$ , we obtain

$$(21) \quad k \int_{a_1}^{\frac{b_2 - b_1}{2t'}} x dx dt' = \frac{k}{2} \left[ \frac{(b_2 - b_1)^2}{4t'^2} - a_1^2 \right] dt'$$

(c) When  $t' \leq -\frac{b_2 - b_1}{2a_2}$ , we similarly obtain

$$(22) \quad k \int_{a_1}^{-\frac{b_2 - b_1}{2t'}} x dx dt' = \frac{k}{2} \left[ \frac{(b_2 - b_1)^2}{4t'^2} - a_1^2 \right] dt'$$

From (18), (19), (20), (21) and (22), the frequency function of  $t'$  is given by

$$(23) \quad F(t') = \frac{a_2 + a_1}{2(b_2 - b_1)} \quad \text{when} \quad -\frac{b_2 - b_1}{2a_2} \leq t' \leq \frac{b_2 - b_1}{2a_2};$$

$$(24) \quad F(t') = \frac{1}{2(b_2 - b_1)(a_2 - a_1)} \left[ \frac{(b_2 - b_1)^2}{4t'^2} - a_1^2 \right],$$

where the range of  $t'$  is subject to either the inequalities,

$$\frac{b_2 - b_1}{2a_2} \leq t' \leq \frac{b_2 - b_1}{2a_1}, \quad \text{or} \quad -\frac{b_2 - b_1}{2a_1} \leq t' \leq -\frac{b_2 - b_1}{2a_2}.$$

See Fig. 4 for the general form of the  $F(t')$  frequency curve. If we make  $a_1 = 0$ , the curve becomes infinite in range. If we make not only  $a_1 = 0$ , but  $(b_1 + b_2)/2 = 0$ , we have, in place of (17),

$$y = mx.$$

In this limiting situation, if we make  $a_2 = a$  and  $\frac{b_2 - b_1}{2} = b$ ,

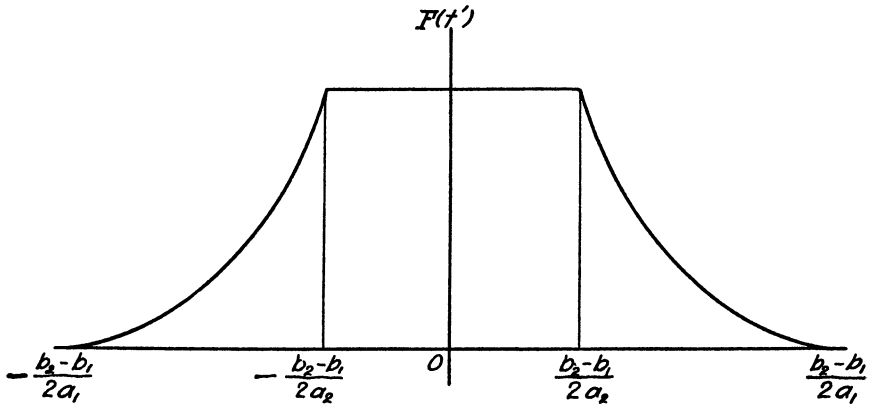


FIG. 4

(23) becomes

$$(25) \quad F(t') = \frac{a}{4b}, \quad \text{for} \quad -\frac{b}{a} \leq t' \leq \frac{b}{a}, \quad \text{and (24) becomes}$$

$$(26) \quad F(t') = \frac{b}{4at'^2} \quad \text{for} \quad t' \geq \frac{b}{a} \quad \text{and for} \quad t' \leq -\frac{b}{a}.$$

Then we have  $y' = y - mx$

and 
$$t' = \frac{y'}{x} = t - m.$$

Further, if  $t'$  is distributed in accord with a frequency function,  $F(t')$ , the distribution of  $t = t' + m$  with  $m$  constant is given by

$$F(t - m).$$



Hence, the probability that a random value  $t$  will fall into a range  $t$  to  $t + dt$  is given to within infinitesimals of higher order by

$$(27) \quad \frac{a}{4b} dt \quad \text{when} \quad m - \frac{b}{a} \leq t \leq m + \frac{b}{a},$$

and by

$$(28) \quad \frac{b dt}{4a(t - m)^2} \quad \text{when} \quad t \geq m + \frac{b}{a} \quad \text{and} \quad t \leq m - \frac{b}{a}.$$

With the frequency curve given by (27) and (28) we may note that the variance of  $t$  becomes infinite.

Without taking the space to continue illustrations, it is fairly obvious that a wide diversity of form can be given to the frequency function of the quotients  $t = y/x$  by relatively simple changes in the location of a sample parent population with reference to the origin.