

ON THE DISTRIBUTION OF THE "STUDENT" RATIO FOR SMALL SAMPLES FROM CERTAIN NON-NORMAL POPULATIONS¹

BY H. L. RIETZ

Much of interest in the theory and practice of statistical methods has been developed around the distribution function,

$$(1) \quad \frac{\Gamma(N/2)}{\pi^{1/2} \Gamma\left(\frac{N-1}{2}\right) (1+z^2)^{N/2}}$$

of the "Student" ratio, $z = \frac{\bar{x} - m}{s}$, where \bar{x} denotes the mean, s the standard deviation of a sample of N items, say x_1, x_2, \dots, x_N , taken at random from a normally distributed parent population of mean, m .

The investigations of certain non-normal parent distributions by Shewhart and Winters [1], Rider [2], E. S. Pearson [3], M. S. Bartlett [4], and R. C. Geary [5] indicate that applications of the "Student" theory give more satisfactory results than the classical theory for a considerable variety of non-normal parent distributions, but some of these investigators find that the theory fails in certain cases to describe the facts to an extent that suggests further experimental sampling investigations along this line whenever suitable data are available. Others infer that a completely satisfactory analysis of the position of the "Student" z -test will be possible only if the theoretical distribution of z in samples from the non-normal distribution in question becomes known. Several of the above named statisticians have attributed the failures of the distribution (1) to describe their data, in large part, to the correlation between $x = \bar{x} - m$ and s . For this reason, there is considerable interest in the degree of correlation between $x = \bar{x} - m$ and s , and especially in the nature of the regression of s or of s^2 on x .

The present paper gives an analysis of data obtained by experimental sampling from two non-normal distributions whose sources we shall now describe. The parent distributions with which the paper is concerned are theoretical distributions resulting from certain urn schemata devised [6] by the writer some years ago.

In 1925, Leone E. Chesire, in an unpublished thesis for the degree, Master of Science, at the University of Iowa, obtained data by experimental sampling, that seem to be appropriate material for a study of the correlation of mean and standard deviation for small samples from certain non-normal distributions.

One of the original bivariate parent populations, whose marginal totals we are

¹ Presented in part before the American Mathematical Society under a somewhat different title, November 26, 1937.

using, exhibited linear regression while the other exhibited non-linear regression. For convenience in distinguishing between the two cases, we shall speak of material from the linear case as Case I and that from the non-linear case as Case II. After devising a scheme for drawing pairs of variates at random, 5,000 pairs were drawn in sets of five for each of the two cases.

While the primary purpose of this experimental sampling was to study the distributions of means, standard deviations, and correlation coefficients [7] for small samples from the non-normal populations, we have as a by-product, in the marginal totals of the correlation tables, four sets of 1,000 pairs of means and standard deviations. However, since three of the four sets of marginal totals of the two theoretical parent correlations tables are alike, we have actually only two significantly different sets to consider.

Case I. For the case of linear regression of y on x in the bivariate parent population, the parent distribution from the marginal totals may be very simply described by showing the frequency distribution in Table 1.

TABLE 1

Sums in second throw of dice-values of stochastic variable.....	2	3	4	5	6	7	8	9	10	11	12
Frequency.....	6	12	18	24	30	36	30	24	18	12	6

The moment coefficients and β 's which characterize the distribution given in Table 1 are:

$$\text{Mean} = 7, \quad \mu_2 = 5\frac{5}{8}, \quad \mu_3 = 0, \quad \mu_4 = 80.5, \quad \beta_1 = 0, \quad \beta_2 = 2\frac{64}{175}.$$

Each of the 1000 sets of five drawn from the distribution in Table 1, yields a mean \bar{y} and a standard deviation, s_y , which we shall denote by w to make our notation simpler to write. Table 2 is the correlation table of the pairs (\bar{y}, w) . The correlation coefficient $r_{w\bar{y}}$, between mean \bar{y} and standard deviation $s_y = w$ has a value

$$r_{w\bar{y}} = -0.020 \pm 0.021$$

which differs insignificantly from zero.

The uncorrected value of the correlation ratio of w on \bar{y} is

$$\eta_{w\bar{y}} = 0.182.$$

When we remember that the correlation ratio is not free to vary in the negative direction from 0, and apply the Pearson correction [8] for this situation together with the "Student" correction [9] for grouping, we obtain for the corrected, $\eta_{w\bar{y}}$, the value 0.133.

It becomes fairly obvious that significant correlation exists and that the regression is non-linear. Indeed, it has been shown recently by Geary [5, pp. 178-9] that normality in the parent distribution is both a necessary and

TABLE 2

Correlation of mean \bar{y} , and standard deviation $s_y = w$, of samples of five items for Case I. Mean of \bar{y} 's = $\bar{y} = 7.141$. Correlation coefficient $r_{w\bar{y}} = -0.020 \pm 0.021$, $s_y = w = 2.079$. Correlation ratio of w on \bar{y} , $\eta_{w\bar{y}} = 0.182$ (uncorrected).

	\bar{y}																	
	3.7	4.1	4.5	4.9	5.3	5.7	6.1	6.5	6.9	7.3	7.7	8.1	8.5	8.9	9.3	9.7	10.1	f_w
4.1								1		2	1							4
3.9							1	1			1	1						4
3.7				1				3		1			2					7
3.5			1					1	1	2	1	1	1					8
3.3					2	4	2	1		7	2	2	5	1				26
3.1				1	3	2	2	8	3	9	4	6	3					41
2.9					5	4	3	10	8	10	3	1	4	2				50
2.7			1	1		3	12	9	18	18	6	9			4			81
2.5				2		4	8	14	15	8	22	14	14	6	5	1		113
2.3					1	2	9	14	10	11	10	10	12	7	3			89
2.1			3	4	2	12	12	16	13	15	8	9	7	8	3			112
1.9				1		7	6	5	16	7	18	15	7	4	3		3	101
1.7					7	7	9	15	15	23	16	14	17	6	4		2	135
1.5				1	3	5	4	3	6	8	6	7	10	6	5	1	1	68
1.3				2	2	1	3	4	11	8	7	11	3	3	2	2		59
1.1					1	1	4	5	5	10	6	6	6	6	1		1	52
0.9				2	1			5	4	2		2	4		1		1	22
0.7						1				7	2	3	2					16
0.5									3	1	1	1		2		2		10
0.3																		1
0.1																		1
$f_{\bar{y}}$	1	1	13	26	40	68	98	135	130	152	109	104	62	35	17	7	2	1000

sufficient condition for the independence of the mean and standard deviation in samples.

Since the number of correlated items, $N = 1000$, is fairly large, we examine into the significance of $\eta_{w\bar{y}} = 0.182$ under the assumption that $N\eta_{w\bar{y}}^2$ is approximately distributed [10] as χ^2 with $a - 1 = 16$ degrees of freedom. This criterion gives odds in favor of significant correlation on approximately a 100 to 1 level of probability.

Next, the means of arrays, \bar{w}_p , were plotted to scale on Table 2 to give a general notion of the nature of the regression of $w = s_y$ on \bar{y} . The location of these means of arrays of w 's affords at least a suggestion of parabolic regression [11] with the curvature concave downward as is to be expected when $\beta_2 - \beta_1 - 3 < 0$, where the β 's relate to the parent distribution.

The next step taken was to analyze the variance, as indicated partly in Table 3, where w_i ($i = 1, 2, \dots, N$) denotes the stochastic variates, a the number of arrays of w 's, \bar{w} the mean of the N values of w_i , n_p ($p = 1, 2, \dots, a$) the number of variates in an array marked p , \bar{w}_p the mean of the array marked p , and where the class interval in Table 2, is taken as the unit.

TABLE 3

	Sum of squares	
For deviations of means of arrays of w 's.....	$\sum_{p=1}^a n_p(\bar{w}_p - \bar{w})^2 = 380$	$a - 1 = 16$
For deviations of variates from the means of their arrays.....	$\sum \sum (w_i - \bar{w}_p)^2 = 11,098$	$N - a = 983$
Total.....	$\sum_{i=1}^N (w_i - \bar{w})^2 = 11,478$	$N - 1 = 999$

In the exhibit given in Table 3, we use the usual algebraic identity

$$(2) \quad \sum_{i=1}^N (w_i - \bar{w})^2 = \sum_{p=1}^a n_p(\bar{w}_p - \bar{w})^2 + \sum \sum (w_i - \bar{w}_p)^2,$$

where the double sum is made up of a sum of N squares.

By dividing the members of (2) by N , we have

$$(3) \quad \frac{1}{N} \sum_{i=1}^N (w_i - \bar{w})^2 = \frac{1}{N} \sum_{p=1}^a n_p(\bar{w}_p - \bar{w})^2 + \frac{1}{N} \sum \sum (w_i - \bar{w}_p)^2$$

The writer has used the identity (3) for many years in lectures to beginners in statistics in proving the equivalence of two definitions of the correlation [12] ratio and is strongly of the opinion that the equality in form (3) appeals more readily to the intuitions of many readers, because of their acquaintance with statements in the language of averages, than does the equivalent equality (2) in the language of sums of squares.

In an extended and more compact form, the analysis is shown in the standard form in Table 4.

TABLE 4

Variance	Degrees of freedom	Sum of squares	Mean square	z -test
Between arrays.....	16	380	23.75	$\frac{1}{2} \log_e 23.75 = 1.584$
Within arrays.....	983	11,098	11.29	$\frac{1}{2} \log_e 11.29 = 1.212$
Total.....	999	11,478		Difference = 0.372

When the sum of squares equal to 380 associated with variance between arrays is further analyzed into a part which could be represented by linear regression,

and a part which represents deviations of the calculated means of arrays of w 's from a straight regression line of w on \bar{y} , the deviations being measured parallel to the w -axis, we find that the part of the amount 380 represented by linear regression is given by

$$Nr_{w\bar{y}}^2 s_w^2 = 1000 (.00040)(11.487) = 4.3.$$

Since both $r = .020 \pm 0.021$ and the small value, 4.3, as part of the sum of squares amounting to 380, may well be regarded as sampling fluctuations, we revert to the figures in Table 3 and apply the Fisher z -test. It turns out that the correlation is significant on practically the 100 to 1 level of probability which conforms well with the above inference based on the assumption that $Nr_{w\bar{y}}^2$ is distributed as χ^2 , with $a - 1$ degrees of freedom.

Next, we computed 1000 values of the "Student" ratio $z = (\bar{y} - 7)/w$, for Case I. One of these 1000 values was of the indeterminate form $\frac{0}{0}$. A frequency distribution of the 999 determinate ratios is shown in column (3), Table 5.

By grouping together the class frequencies at the tails of the theoretical distributions until each of the end class frequencies is not less than 5, and calculating χ^2 for the observed distribution in column (3) in comparison with the theoretical distribution in column (6) as found from the "Student" theory in samples of 5 items from a normal distribution, we obtain $\chi^2 = 3.728$ with 11 degrees of freedom.

Thus, the differences between the distribution in column (3) and the "Student" distribution for $N = 5$ shown in column (6) are not only insignificant under the χ^2 -test, but are so small as to be expected in a relatively small percentage of statistical experiments even if the "Student" z -distribution were the theoretically exact distribution of our ratios.

The usual moment coefficients of the distribution of observed z 's in column (3), Table 5, are:

$$\begin{aligned} \mu'_1 &= 0.033533, & \mu_3 &= 0.254383, & \beta_1 &= 0.55955, \\ s = \sqrt{\mu_2} &= 0.69799, & \mu_4 &= 2.22504, & \beta_2 &= 9.37353. \end{aligned}$$

Since the value, 0.69799, of the standard deviation of the observed distribution differs very little from $1/\sqrt{N-3} = 0.70711$, the normal curve fitted by using the standard deviation of the observed distribution (column 4, Table 5) differs very little from the normal curve with the origin at the population mean and standard deviation, $\sqrt{2}/2$, (column 5). Furthermore, the application of the χ^2 -test to columns (4) and (5) of Table 5 with class frequencies in the "tails" grouped as above gives $\chi^2 = 2.91$ with 9 degrees of freedom.

The moment coefficients of the observed distribution indicate a markedly leptokurtic and somewhat skew distribution but the indications of skewness may be traced mainly and perhaps entirely to the presence of the two extreme variates at the upper end of the distribution and separated about three times the standard deviation from the next class frequency that differs from zero. By

TABLE 5
Distribution of the ratios, $z = (\bar{y} - 7)/w$ in samples of $N = 5$ for Case I.

(1)	(2)	(3)	(4)	(5)	(6)
$z = (\bar{y} - 7)/w$	$t = z\sqrt{N-1} = 2z$	Observed distribution	Normal distribution fitted to observed column (3)	Normal distribution of S.D. $= \frac{1}{\sqrt{N-3}} = \frac{1}{\sqrt{2}}$ in same units as z (measured from population mean)	From the Student theoretical distribution for $N = 5$
-6.0	-12				0.1
-5.5	-11				0.1
-5.0	-10				0.1
-4.5	-9				0.2
-4.0	-8				0.3
-3.5	-7				0.6
-3.0	-6	2	0.05	0.1	1.3
-2.5	-5	1	0.75	0.4	2.7
-2.0	-4	5	3.6	6.2	7.0
-1.5	-3	17	27.7	32.0	21.0
-1.0	-2	67	98.5	105.9	70.5
-0.5	-1	216	210.8	217.2	217.5
0	0	357	279.3	275.4	356.2
0.5	1	226	225.7	217.2	217.5
1.0	2	75	111.7	105.9	70.5
1.5	3	22	33.7	32.0	21.0
2.0	4	5	6.2	6.2	7.0
2.5	5	1	0.75	0.4	2.7
3.0	6	3	0.05	0.1	1.3
3.5	7	0			0.6
4.0	8	0			0.3
4.5	9	0			0.2
5.0	10	2			0.1
5.5	11				0.1
6.0	12				0.1
		999	998.8	999.0	999.0

excluding these two variates from our calculations, we obtain the following moment coefficients:

$$\mu'_1 = 0.023571, \quad \mu_3 = 0.022264, \quad \beta_1 = 0.0058786,$$

$$s = \sqrt{\mu_2} = 0.662202, \quad \mu_4 = 1.009673. \quad \beta_2 = 5.2507062.$$

In the observed distribution thus modified, by excluding the extreme upper class frequency 2, the evidence of skewness has disappeared.

Case II. For our Case II we have a frequency distribution as shown in Table 6.

TABLE 6

Totals in second throws of two dice-values of the stochastic variable....	2	3	4	5	6	7	8	9	10	11	12
Frequency	1	4	9	16	25	36	35	32	27	20	11

Again, since with the uncorrected $\eta_{v\bar{u}}$, Table 6, we have $N\eta_{v\bar{u}}^2 = 31.5$, and since $N\eta_{v\bar{u}}^2$ is approximately distributed as χ^2 with $a - 1 = 17$ degree of freedom, we have odds of the order of 100 to 1 against so large a value being a mere sampling fluctuation.

TABLE 7'

Correlation of mean \bar{u} , and standard deviation $s_u = v$, of five items for Case II, mean of $\bar{u} = \bar{\bar{u}} = 6.971$, Correlation coefficient $r_{v\bar{u}} = -0.012 \pm 0.020$.
 $v = s_u = 2.044$. Correlation ratio of v on \bar{u} , $\eta_{v\bar{u}} = 0.177$ (uncorrected).

v

	3.7	4.1	4.5	4.9	5.3	5.7	6.1	6.5	6.9	7.3	7.7	8.1	8.5	8.9	9.3	9.7	10.1	10.5	10.9	f_v
3.9								1		1	1									3
3.7							1				3	1								5
3.5					1	2	1	4	4	2		2								16
3.3						1	5	4	3	5	2	2	1	1	1					25
3.1						4	6	9	6	4	6	1	1	1						38
2.9					3		8	10	16	8	8	5	5		2					65
2.7				1	2	4	10	7	7	17	13	4	1	2	1					69
2.5			1	1	2	10	10	17	8	19	11	5	4	3						91
2.3		1		3	5	5	11	11	14	13	10	7	10	3		1				94
2.1		1	4		3	15	21	16	13	15	13	9	5	3						118
1.9	1		1	1	12	12	7	12	9	16	19	7	13	5						115
1.7		1	1	6	5	11	14	12	16	17	20	10	7	5	2	3				130
1.5				1	7		14	8	8	3	6	5	4	3	3	1	1			64
1.3			1	1	7	10	3	11	3	13	4	1	6	1	1	2				64
1.1		1		2	3	4	5	7	9	9	4	4	5		3	1			1	58
0.9					2	1	2	4	4	6	2	3	1	3	1	1				30
0.7						2	3		1		2	2								10
0.5								2	1	2										5
Σ		4	8	16	52	81	121	135	122	150	124	68	63	30	14	9	-	0	-	1000

Now proceeding to the analysis of variance, we substitute our numerical values derived from Table 7 in the identity

$$(4) \quad \sum_{i=1}^N (v_i - \bar{v})^2 = \sum_{p=1}^a n_p (\bar{v}_p - \bar{v})^2 + \sum \sum (v_i - \bar{v}_p)^2$$

and obtain, in terms of class intervals as units,

$$10,871 = 340 + 10,531.$$

An outline of the analysis is exhibited in Table 8

TABLE 8

Variance	Degrees of freedom	Sum of squares	Mean square	<i>z</i> -test
Between arrays.....	17	340	20.00	$\frac{1}{2} \log_e 20.00 = 1.50$
Within arrays.....	982	10,531	10.72	$\frac{1}{2} \log_e 10.72 = 1.18$
Total.....	999	10,871		Diff. = 0.32

The moment coefficients and β 's which characterize the distribution in Table 6 are:

$$\begin{aligned} \text{Mean} &= 7.972, & \mu_2 &= 4.888, & \mu_3 &= -1.755, & \mu_4 &= 58.724, \\ & & \beta_1 &= 0.0264, & \beta_2 &= 2.449. \end{aligned}$$

As in the linear case, samples of 5000 pairs of variates were drawn in sets of five by Miss Chesire. Analogous to Case I, our first concern is with the regression of the standard deviation, $s_u = v$, of u from a sample of five on its mean value, \bar{u} .

The correlation table for values of \bar{u} and v is shown in Table 7. The correlation coefficient is $r_{v\bar{u}} = -0.012 = \pm 0.021$, but the uncorrected correlation ratio of v on \bar{u} is given by

$$\eta_{v\bar{u}} = 0.177.$$

After applying the Pearson and Student corrections, we obtain the corrected

$$\eta_{v\bar{u}} = 0.131.$$

When the sum of squares, 340, associated with variance between arrays is further analyzed into a part which could be represented by linear regression, and a part which represents deviations of the calculated means of arrays of v 's from a straight regression line of v on \bar{u} , the deviations being measured parallel to the v -axis, we find that the part of the amount 340 represented by linear regression, would be only $Nr^2s_v^2 = 1000 (.000144)(10.871) = 1.6$.

Since both $r_{v\bar{u}} = -0.012 \pm 0.021$ and the small value, 1.6, as part of the sum of squares 340, may well be regarded as sampling fluctuations, we revert to the figures of Table 8.

The difference of the logarithms in the last column of Table 8, is 0.32, which corresponds to a level of significance of the general order of 100 to 1. Next, we calculate and plot on Table 7 the means of arrays of v 's to give a general notion of the regression of v on \bar{u} . The location of these means of arrays suggests rather strongly that the regression of v on \bar{u} is parabolic with the curvature concave downward as we should expect from the fact that $\beta_2 - \beta_1 - 3 < 0$, where the β 's pertain to the parent distribution.

Next, we computed 1000 values of the "Student" ratio, $z = (u - 7.972)/v$,

for Case II. One of these ratios was infinite. A frequency distribution of the 999 determinate ratios is shown in column 3, Table 9.

The observed distribution (column 3) and the "Student" distribution (column 6) of Table 9, to be expected in samples of $N = 5$, when samples are drawn from a normal distribution, are in close agreement. In fact, when we group together the tail frequencies of the theoretical distribution until each of them is not less than 5, the result of testing the goodness of fit gives $\chi^2 = 17.187$ with 11 degrees of freedom. This gives a value in the neighborhood of 0.1 for the probability, P , that as large or larger deviations than that experienced will occur, due to chance fluctuations, in a single repetition of the experiment. In other words, on the basis of this test, the indications are that we should have in the long run, as large or larger deviations than we have experienced in this case, in about 10 per cent of a large number of sets of sampling of 1000 per set even when the sampling is from a normal distribution.

TABLE 9
Distribution of the ratio, $(\bar{u} - 7.972)/v$ in samples of five for Case II.

(1) $z = (\bar{u} - 7.972)/v$	(2) $t = \frac{z\sqrt{N-1}}{2z}$	(3) Observed	(4) Normal distribution fitted to observed, Column (3).	(5) Normal distribution with S.D. = $\frac{1}{\sqrt{N-3}}$ and origin at population mean	(6) Student's z -distribution for normal parent population with $N = 5$
-5.5	-11	1			0.1
-5.0	-10				0.1
-4.5	-9				0.2
-4.0	-8				0.3
-3.5	-7				0.6
-3.0	-6		0.1	0.1	1.3
-2.5	-5	2	0.4	0.4	2.7
-2.0	-4	3	4.3	6.2	7.0
-1.5	-3	23	25.4	32.0	21.0
-1.0	-2	48	92.0	105.9	70.5
-0.5	-1	203	205.3	217.2	217.5
-0.0	0	380	278.4	275.4	356.2
0.5	1	226	231.4	217.2	217.5
1.0	2	72	117.5	105.9	70.5
1.5	3	24	36.5	32.0	21.0
2.0	4	9	6.9	6.2	7.0
2.5	5	3	0.8	0.4	2.7
3.0	6	4	0.1	0.1	1.3
3.5	7	1			.6
4.0	8				.3
4.5	9				.2
					.1
					.1
					.1
Total		999	999.1	999.0	999.0
∞	∞	1			

SUMMARY

1. The linear correlation coefficient, r , of the mean and standard deviation differs insignificantly from 0 in each case.

2. The correlation ratio of the standard deviation on the mean differs significantly from 0, and the regression of the standard deviation on the mean conforms, in its general aspects, to expectation under the theory of Neyman [12].

3. The indeterminate "Student" ratio of the form, $\frac{0}{0}$, in Case I and that of the form, (constant)/0, in Case II are probably due in part to grouping into class intervals, but the infinite ratio would undoubtedly have had such a large value that it would be excluded from calculations under any one of the known criteria for rejection of extreme observations.

4. Although the rejection of one indeterminate ratio in each of the two cases is slightly disturbing, the evidence presented by our analysis of the experimental sampling lends support to the view that the results of the "Student" theory are almost certainly applicable, for many purposes, when the parent distributions are of such non-normal types as are involved in our sampling.

REFERENCES

- [1] W. A. Shewhart and F. W. Winters, "Small samples—new experimental results," *Journal American Statistical Association*, Vol. 23 (1928), pp. 144–153.
- [2] P. R. Rider, "On the distribution of the ratio of mean to standard deviation in small samples from non-normal universes," *Biometrika*, Vol. 21 (1929), pp. 124–143. "On small samples from certain non-normal universes," *Annals of Mathematical Statistics*, Vol. 2, (1931), pp. 48–65.
- [3] E. S. Pearson, "The distribution of frequency constants in small samples from non-normal symmetrical and skew populations," *Biometrika*, Vol. 21 (1929), pp. 259–286.
- [4] M. S. Bartlett, "The effect of non-normality on the t distribution," *Proceedings of the Cambridge Philosophical Society*, Vol. 31 (1935), pp. 223–231.
- [5] R. C. Geary, "The distribution of "Student's" ratio for non-normal samples," *Supplement to the Journal of the Royal Statistical Society*, No. 2, 1936, pp. 178–184.
- [6] H. L. Rietz, "Urn schemata as a basis for the development of correlation theory," *Annals of Mathematics*, Vol. 21 (1920), pp. 306–322.
- [7] E. S. Pearson, Leona Chesire, and Elena M. Oldis, "Further experiments on the sampling distribution of the correlation coefficient," *Journal American Statistical Association*, Vol. 27 (1932), pp. 121–128.
- [8] Karl Pearson, "On a correction to be made in the correlation ratio," *Biometrika*, Vol. 8, (1911–12), pp. 254–6.
- [9] "Student," "The correction to be made in the correlation ratio for grouping," *Biometrika*, (1913), pp. 316–20.
- [10] R. A. Fisher, *Statistical methods for research workers*, Fourth Edition, p. 237.
- [11] J. Neyman, "On the correlation of mean and variance in samples drawn from an 'infinite' population," *Biometrika*, Vol. 18 (1926), pp. 401–13.
- [12] H. L. Rietz, *Mathematical Statistics* (Carus Monograph), 1926, p. 91.