

ON THE PROBABLE ERROR OF A COEFFICIENT OF CORRELATION AS FOUND FROM A FOURFOLD TABLE.

By KARL PEARSON, F.R.S.

LET the fourfold table be

a	b	$a+b$
c	d	$c+d$
$a+c$	$b+d$	N

Then on the assumption that the frequency distribution is normal, we can by aid of Everitt's Tables of the Tetrachoric Functions* rapidly find r . I have shown in a paper published in the *Phil. Trans.* in 1900† that found in this way

Probable error of r

$$= \frac{.67449}{\sqrt{N}\chi_0} \left\{ \frac{(a+d)(c+b)}{4N^2} + \psi_2^2 \frac{(a+c)(d+b)}{N^2} + \psi_1^2 \frac{(a+b)(d+c)}{N^2} \right. \\ \left. + 2\psi_1\psi_2 \frac{ad-bc}{N^2} - \psi_2 \frac{ab-cd}{N^2} - \psi_1 \frac{ac-bd}{N^2} \right\}^{\frac{1}{2}} \dots(i),$$

where

$$\psi_1 = \frac{1}{\sqrt{2\pi}} \int_0^{\beta_1} e^{-\frac{1}{2}z^2} dz, \quad \psi_2 = \frac{1}{\sqrt{2\pi}} \int_0^{\beta_2} e^{-\frac{1}{2}z^2} dz,$$

$$\beta_1 = \frac{h-rk}{\sqrt{1-r^2}}, \quad \beta_2 = \frac{k-rh}{\sqrt{1-r^2}},$$

$$\chi_0 = \frac{1}{2\pi} \frac{1}{\sqrt{1-r^2}} e^{-\frac{1}{2} \frac{1}{1-r^2} (h^2 + k^2 - 2rhk)},$$

* *Biometrika*, Vol. VII, p. 436, and Vol. VIII, p. 385.

† *Phil. Trans.* Vol. 195 A, p. 14. Owing to the carelessness of the printers my χ_0 was put as $\sqrt{\chi_0}$ and the last N^2 in the denominator as N_2 .

and h and k have their usual meaning defined by the integrals

$$\frac{(a+c)-(b+d)}{2N} = \frac{1}{\sqrt{2\pi}} \int_0^h e^{-\frac{1}{2}z^2} dz = \frac{1}{2}\alpha_1, \text{ say;}$$

$$\frac{(a+b)-(c+d)}{2N} = \frac{1}{\sqrt{2\pi}} \int_0^k e^{-\frac{1}{2}z^2} dz = \frac{1}{2}\alpha_2, \text{ say.}$$

Let $H = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}h^2}$, $K = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}k^2}$ as usual.

Now the formula (i) above for the probable error of r is admittedly laborious in use. I have tried in many ways, while retaining its full accuracy, to throw it into a form involving less laborious calculations; I have not succeeded, however, in achieving any sensible reduction in its complexity, as long as I maintain its complete generality.

Although many hundred fourfold tables have now been published, many of which give such small correlations that their true significance can only be settled by a knowledge of their probable errors, yet I find only 40 to 50 probable errors have so far been determined. This matter seems so regrettable that I have sought for a fairly easy method of determining a closely empirical expression for the probable error of r which is likely to be of service, and can be adapted easily to tables.

I consider first two extreme cases. If h and k are both zero, or the fourfold division at the mean, then $\psi_1 = \psi_2 = 0^*$,

Probable error of r

$$= \frac{.67449 \cdot 2\pi \sqrt{1-r^2} \left\{ \frac{(a+d)(b+c)}{4N^2} \right\}^{\frac{1}{2}}}{\sqrt{N}} = \frac{.67449 \sqrt{1-r^2} \pi}{\sqrt{N}} \frac{1}{2} \left\{ \frac{16ab}{N^2} \right\}^{\frac{1}{2}},$$

since in this case $a=d$, and $b=c$.

But for a division at the mean by Sheppard's Theorem

$$r = \cos \pi \frac{b}{a+b} = \sin \left(\frac{\pi}{2} - \frac{\pi b}{a+b} \right),$$

or

$$(\sin^{-1} r) / \frac{1}{2}\pi = (a-b)/(a+b).$$

Hence

$$1 - \left(\frac{\sin^{-1} r}{\frac{1}{2}\pi} \right)^2 = \frac{4ab}{(a+b)^2} = \frac{16ab}{N^2}.$$

Substituting we have:

$$\text{Probable error of } r = \frac{.67449 \pi}{\sqrt{N}} \frac{1}{2} \sqrt{1-r^2} \sqrt{1 - \left(\frac{\sin^{-1} r}{90^\circ} \right)^2} \dots\dots(ii),$$

if the angle of the inverse sine be read in degrees.

Again if $r=0$, the probable error of r may be obtained from (i) whatever the values of h and k . For in this case

$$ad-bc=0, \quad \psi_1 = \frac{1}{2}\alpha_1, \quad \psi_2 = \frac{1}{2}\alpha_2, \quad \chi_0 = HK.$$

* *Phil. Trans.* Vol. 192 A, p. 141 and Vol. 195 A, p. 7.

24 On the Probable Error of a Coefficient of Correlation

We have $(b+d)/N = \frac{1}{2}(1 - \alpha_1)$, $(a+c)/N = \frac{1}{2}(1 + \alpha_1)$,

$$(a+b)/N = \frac{1}{2}(1 + \alpha_2), \quad (c+d)/N = \frac{1}{2}(1 - \alpha_2),$$

$$\begin{aligned} \frac{a+d}{N} &= \frac{(a+b)(a+c)}{N^2} + \frac{ad-bc}{N^2} + \frac{(c+d)(b+d)}{N^2} + \frac{ad-bc}{N^2} \\ &= \frac{1}{4}(1 + \alpha_2)(1 + \alpha_1) + \frac{1}{4}(1 - \alpha_2)(1 - \alpha_1) = \frac{1}{2}(1 + \alpha_1\alpha_2), \end{aligned}$$

since $ad-bc=0$ in the original population.

$$\text{Similarly:} \quad \frac{c+b}{N} = \frac{1}{2}(1 - \alpha_1\alpha_2).$$

$$\begin{aligned} \frac{ab-cd}{N^2} &= \frac{a(N-a-c-d)-cd}{N^2} \\ &= \frac{a}{N} - \frac{(a+c)(a+d)}{N^2} \\ &= \frac{1}{4}(1 + \alpha_2)(1 + \alpha_1) - \frac{1}{4}(1 + \alpha_1)(1 + \alpha_1\alpha_2) \\ &= \frac{1}{4}\alpha_2(1 - \alpha_1^2), \end{aligned}$$

$$\text{and similarly:} \quad \frac{ac-bd}{N^2} = \frac{1}{4}\alpha_1(1 - \alpha_2^2).$$

Hence substituting in (i)

$$\begin{aligned} \text{Probable error of } r &= \frac{.67449}{\sqrt{N}HK} \left\{ \frac{1}{16}(1 - \alpha_1^2\alpha_2^2) + \frac{1}{16}\alpha_2^2(1 - \alpha_1^2) + \frac{1}{16}\alpha_1^2(1 - \alpha_2^2) \right. \\ &\quad \left. - \frac{1}{8}\alpha_2^2(1 - \alpha_1^2) - \frac{1}{8}\alpha_1^2(1 - \alpha_2^2) \right\}^{\frac{1}{2}} \\ &= \frac{.67449}{\sqrt{N}HK} \sqrt{\frac{1}{16}(1 - \alpha_1^2)(1 - \alpha_2^2)} \dots\dots\dots(\text{iii}). \end{aligned}$$

This can also be put in the form:

$$\text{Probable error of } r = \frac{.67449}{\sqrt{N}HK} \sqrt{\frac{(a+b)(a+c)(d+b)(d+c)}{N^4}} \dots\dots(\text{iv}).$$

This is the probable error of r of a fourfold table when the real value of r is zero.

Now as (ii) and (iv) give the reducing factors for the two cases (a) when h and k are both zero but r has any value and (b) when h and k have any values but r is zero, it occurred to me that the combined product of the two would give good results for a considerable range of values of h and k and r . We have to note that (iv) for h and k zero becomes

$$\frac{.67449}{\sqrt{N}} \frac{\pi}{2}.$$

Hence we take as our formula:

Probable error of r

$$= \frac{.67449}{\sqrt{N}} \sqrt{1 - r^2} \sqrt{1 - \left(\frac{\sin^{-1} r}{90^\circ} \right)^2} \sqrt{\frac{1}{2}(1 + \alpha_1) \frac{1}{2}(1 - \alpha_1)} \sqrt{\frac{1}{2}(1 + \alpha_2) \frac{1}{2}(1 - \alpha_2)} \dots(\text{v}).$$

Now it will be seen that this consists of three parts:

- (a) $\sqrt{1-r^2} \sqrt{1 - \left(\frac{\sin^{-1} r}{90^\circ}\right)^2}$. This is easy to table for all values of r .
- (b) $\frac{\sqrt{\frac{1}{2}(1+\alpha_1)\frac{1}{2}(1-\alpha_1)}}{H}$, and
- (c) $\frac{\sqrt{\frac{1}{2}(1+\alpha_2)\frac{1}{2}(1-\alpha_2)}}{K}$.

Both these (b) and (c) can be readily found from a single table rapidly formed from Sheppard's Table of the Probability Integral. The entry to the single table will be $(a+c)/N$ or $(a+b)/N$, i.e. $\frac{1}{2}(1+\alpha)$.

Thus a knowledge of the correlation r and the two division percentages (together with Miss Gibson's Table for $\cdot67449/\sqrt{N}$), will enable us by the aid of the two new tables to rapidly write down four factors whose product gives the required probable error. I have tested the form (v) against the true probable error as found from (i). In all cases it gave results differing only from the true value at most by about one or two units in the third place of figures—a result amply accurate for all practical purposes.

Illustration I.

211·25	153·75	365
152·75	560·25	713
364	714	1078

The correlation was found to be $\cdot5557 \pm \cdot0261$; the probable error from the short formula was $\cdot0265$.

Illustration II.

1562	42	1604
383	94	477
1945	136	2081

The correlation was found to be $\cdot5954 \pm \cdot0272$; the probable error from the short formula was $\cdot0293$.

Illustration III.

455	622	1077
599	1324	1923
1054	1946	3000

The correlation was found to be $\cdot1811 \pm \cdot0210$; the probable error from the short formula was $\cdot0199$.

Illustration IV.

849	665	1514
205	1281	1486
1054	1946	3000

26 *On the Probable Error of a Coefficient of Correlation*

The correlation was found to be $\cdot6633 \pm \cdot0132$; the probable error from the short formula was $\cdot0132$.

Illustration V.

1196	223	1419
318	1263	1581
1514	1486	3000

The correlation was found to be $\cdot8464 \pm \cdot0079$; the probable error from the short formula was $\cdot0079$.

These examples will suffice, I think, to give confidence in the formula and in the tables accompanying this paper. The absence of probable errors from the expressions for fourfold table correlations can no longer be justified on the ground of their great laboriousness.

The following Tables have been calculated by Miss Julia Bell.

Let $\chi_1 = \cdot67449/\sqrt{N}$. This is given by Miss Gibson's Tables, *Biometrika*, Vol. III. p. 387. Let

$$\chi_r = \sqrt{1-r^2} \sqrt{1 - \left(\frac{\sin^{-1} r}{90^\circ} \right)^2},$$

and

$$\chi_a = \frac{1}{H} \sqrt{\frac{1}{2}(1+\alpha) \times \frac{1}{2}(1-\alpha)}.$$

Then

Probable error of $r = \chi_1 \cdot \chi_r \cdot \chi_{a_1} \cdot \chi_{a_2}$.

TABLE I. *Values of χ_r for Values of r .*

r	χ_r	r	χ_r	r	χ_r	r	χ_r	r	χ_r
$\cdot00$	1.0000	$\cdot20$	$\cdot9717$	$\cdot40$	$\cdot8845$	$\cdot60$	$\cdot7298$	$\cdot80$	$\cdot4843$
$\cdot01$	$\cdot9999$	$\cdot21$	$\cdot9688$	$\cdot41$	$\cdot8785$	$\cdot61$	$\cdot7200$	$\cdot81$	$\cdot4687$
$\cdot02$	$\cdot9997$	$\cdot22$	$\cdot9657$	$\cdot42$	$\cdot8723$	$\cdot62$	$\cdot7099$	$\cdot82$	$\cdot4526$
$\cdot03$	$\cdot9994$	$\cdot23$	$\cdot9625$	$\cdot43$	$\cdot8659$	$\cdot63$	$\cdot6997$	$\cdot83$	$\cdot4362$
$\cdot04$	$\cdot9989$	$\cdot24$	$\cdot9591$	$\cdot44$	$\cdot8594$	$\cdot64$	$\cdot6892$	$\cdot84$	$\cdot4192$
$\cdot05$	$\cdot9982$	$\cdot25$	$\cdot9556$	$\cdot45$	$\cdot8527$	$\cdot65$	$\cdot6785$	$\cdot85$	$\cdot4018$
$\cdot06$	$\cdot9975$	$\cdot26$	$\cdot9520$	$\cdot46$	$\cdot8458$	$\cdot66$	$\cdot6675$	$\cdot86$	$\cdot3838$
$\cdot07$	$\cdot9966$	$\cdot27$	$\cdot9482$	$\cdot47$	$\cdot8388$	$\cdot67$	$\cdot6563$	$\cdot87$	$\cdot3652$
$\cdot08$	$\cdot9955$	$\cdot28$	$\cdot9442$	$\cdot48$	$\cdot8315$	$\cdot68$	$\cdot6448$	$\cdot88$	$\cdot3461$
$\cdot09$	$\cdot9943$	$\cdot29$	$\cdot9401$	$\cdot49$	$\cdot8241$	$\cdot69$	$\cdot6331$	$\cdot89$	$\cdot3262$
$\cdot10$	$\cdot9930$	$\cdot30$	$\cdot9358$	$\cdot50$	$\cdot8165$	$\cdot70$	$\cdot6211$	$\cdot90$	$\cdot3057$
$\cdot11$	$\cdot9915$	$\cdot31$	$\cdot9314$	$\cdot51$	$\cdot8087$	$\cdot71$	$\cdot6088$	$\cdot91$	$\cdot2843$
$\cdot12$	$\cdot9899$	$\cdot32$	$\cdot9268$	$\cdot52$	$\cdot8007$	$\cdot72$	$\cdot5962$	$\cdot92$	$\cdot2620$
$\cdot13$	$\cdot9881$	$\cdot33$	$\cdot9221$	$\cdot53$	$\cdot7926$	$\cdot73$	$\cdot5834$	$\cdot93$	$\cdot2387$
$\cdot14$	$\cdot9862$	$\cdot34$	$\cdot9172$	$\cdot54$	$\cdot7842$	$\cdot74$	$\cdot5702$	$\cdot94$	$\cdot2142$
$\cdot15$	$\cdot9841$	$\cdot35$	$\cdot9122$	$\cdot55$	$\cdot7756$	$\cdot75$	$\cdot5568$	$\cdot95$	$\cdot1882$
$\cdot16$	$\cdot9819$	$\cdot36$	$\cdot9070$	$\cdot56$	$\cdot7669$	$\cdot76$	$\cdot5430$	$\cdot96$	$\cdot1605$
$\cdot17$	$\cdot9796$	$\cdot37$	$\cdot9016$	$\cdot57$	$\cdot7579$	$\cdot77$	$\cdot5288$	$\cdot97$	$\cdot1305$
$\cdot18$	$\cdot9771$	$\cdot38$	$\cdot8961$	$\cdot58$	$\cdot7488$	$\cdot78$	$\cdot5144$	$\cdot98$	$\cdot0972$
$\cdot19$	$\cdot9745$	$\cdot39$	$\cdot8904$	$\cdot59$	$\cdot7394$	$\cdot79$	$\cdot4995$	$\cdot99$	$\cdot0585$
								1.00	$\cdot0000$

TABLE II.

Values of χ_α for Values of $\frac{1}{2}(1+\alpha)$.

$\frac{1}{2}(1+\alpha)$	χ_α	$\frac{1}{2}(1+\alpha)$	χ_α	$\frac{1}{2}(1+\alpha)$	χ_α	$\frac{1}{2}(1+\alpha)$	χ_α
<i>50</i>	1·2533	<i>65</i>	1·2877	<i>80</i>	1·4288	<i>95</i>	2·1132
<i>51</i>	1·2535	<i>66</i>	1·2928	<i>81</i>	1·4457	<i>96</i>	2·2740
<i>52</i>	1·2539	<i>67</i>	1·2984	<i>82</i>	1·4641	<i>97</i>	2·5071
<i>53</i>	1·2546	<i>68</i>	1·3044	<i>83</i>	1·4844	<i>98</i>	2·8915
<i>54</i>	1·2556	<i>69</i>	1·3109	<i>84</i>	1·5067	<i>985</i>	3·2097
<i>55</i>	1·2569	<i>70</i>	1·3180	<i>85</i>	1·5315	<i>990</i>	3·7333
<i>56</i>	1·2585	<i>71</i>	1·3256	<i>86</i>	1·5590	<i>991</i>	3·8854
<i>57</i>	1·2604	<i>72</i>	1·3338	<i>87</i>	1·5897	<i>992</i>	4·0639
<i>58</i>	1·2626	<i>73</i>	1·3427	<i>88</i>	1·6245	<i>993</i>	4·2784
<i>59</i>	1·2652	<i>74</i>	1·3523	<i>89</i>	1·6640	<i>994</i>	4·5419
<i>60</i>	1·2680	<i>75</i>	1·3626	<i>90</i>	1·7094	<i>995</i>	4·8779
<i>61</i>	1·2712	<i>76</i>	1·3738	<i>91</i>	1·7623	<i>996</i>	5·3278
<i>62</i>	1·2748	<i>77</i>	1·3859	<i>92</i>	1·8249	<i>997</i>	5·9776
<i>63</i>	1·2787	<i>78</i>	1·3990	<i>93</i>	1·9003	<i>998</i>	7·0465
<i>64</i>	1·2830	<i>79</i>	1·4133	<i>94</i>	1·9937	<i>999</i>	9·3870