

A FORTRAN procedure for Fisher's exact probability test

JAMES E. HAYS

Virginia Mason Research Center
1000 Seneca Street, Seattle, Washington 98101

There are several statistical methods available for comparing two percentages to see if a significant difference exists. One method, which is not generally used to full advantage, is Fisher's exact method (Fisher, 1950). This procedure is commonly applied to 2 by 2 contingency tables where some of the observed frequencies are extremely small and result in correspondingly small expected frequencies. In such cases it is not generally considered that the more commonly applied chi-square approximation to the multinomial has a sufficient degree of accuracy, even when corrected for continuity (Yates, 1934). Cochran (1950), for example, suggests that if any of the expected values are less than five, then the approximation may be poor. Even for such cases where the observed frequencies are small, hand computation of the Fisher exact method is extremely laborious. Application of the procedure to tables with larger observed frequencies results in computations that take a prohibitive amount of time. Consequently, for these cases the chi-square approximation is generally used, even though determination of the exact probabilities would be more desirable. With access to a digital computer, Fisher's method can be employed more often; in fact, it can be employed in all cases where a 2 by 2 contingency table is to be analyzed. The program described here requires relatively few storage locations and, therefore, can be implemented on small-scale computing installations.

To illustrate the use of Fisher's method, assume that the two percentages to be compared are in the fractional form $P_1 = A/(A + C)$ and $P_2 = B/(B + D)$ where $A = \text{MIN}(A, B, C, D)$. Here A, C and B, D are the success/failure dichotomies observed in the two populations which we wish to compare. This may be more easily illustrated in tabular form:

	Pop. 1	Pop. 2
Success	A	B
Failure	C	D

The exact probability of observing two fractions, P_1 and P_2 , when there is no class difference is given by

$$P_0 = \frac{(A + B)! (C + D)! (A + C)! (B + D)!}{A! B! C! D! N!}$$

where $N = A + B + C + D$. To obtain the final probability to use in judging whether there is a significant difference, we must add to the above value the probabilities of more divergent fractions than those observed. The next more divergent situation, assuming that $P_1 < P_2$, is obtained by decreasing A and D and increasing B and C by unity. In general, then, the following equation is used:

$$P_i = \left[\frac{(A + B)! (C + D)! (A + C)! (B + D)!}{N!} \right] \times \left[\frac{1}{(A - i)! (B + i)! (C + i)! (D - i)!} \right]$$

where $i = 0, 1, \dots, k$ and $k =$ the number of more extreme distributions.

Finally, the probability of observing two fractions as much or more divergent than P_1 and P_2 , when there is no difference in the sample populations, is given by:

$$\alpha = \sum_{i=0}^k P_i$$

For the purposes of computation, the general equations shown above are represented as follows:

$$\begin{aligned} \log(p_i) = & \left[\sum_{j1=0}^{A+B} \log(j1) + \sum_{j2=0}^{C+D} \log(j2) + \sum_{j3=0}^{A+C} \log(j3) \right. \\ & + \sum_{j4=0}^{B+D} \log(j4) - \sum_{m=0}^N \log(m) \left. \right] \\ & - \left[\sum_{m1=0}^{A-i} \log(m1) + \sum_{m2=0}^{B+i} \log(m2) \right. \\ & + \sum_{m3=0}^{C+i} \log(m3) + \sum_{m4=0}^{D-i} \log(m4) \left. \right] \\ \alpha = & \sum_{i=0}^k e^{\log(p_i)} \end{aligned}$$

In this way, the effect of rounding errors is minimized.

The program can handle 2 by 2 contingency tables with any observed frequencies. The computation time involved is a function of the magnitude of the observed frequencies, smaller frequencies requiring less computer time.

Input. The program is in the form of a FORTRAN subroutine which receives the observed cell frequencies via a standard FORTRAN call statement from a main program. This subprogram can therefore be combined with any main-line program the user desires. For ease of use in the case where tables are to be read from a card reader, a short main-line program is provided.

Output. The printed output consists of the following: (1) a listing of the input contingency table, (2) the exact probability associated with the input table, (3) the exact probabilities associated with each of the more extreme possible tables, and (4) the total probability of observing two fractions as much or more divergent than those observed in the input table.

The main program, supplied with the exact probability subroutine, will process input data until a blank card is reached. The above-described output is produced for each input case provided.

Computer and Language. The program was written in FORTRAN and has been tested on a Raytheon 704.

Availability. Program, test data, and sample output listings and instructions may be obtained free of charge from James E. Hays, Virginia Mason Research Center, 1000 Seneca Street, Seattle, Washington 98101.

REFERENCES

- FISHER, R. A. *Statistical methods for research workers* (11th ed., revised). New York: Hafner Publishing Company, 1950.
- COCHRAN, W. G. The comparison of percentages in matched samples. *Biometrika*, 1950, **37**, 256-266.
- YATES, F. Contingency tables involving small numbers and the Chi Square test. *Supplement to the Journal of the Royal Statistical Society*, 1934, **1**, 217-235.