

Algorithm 434

Exact Probabilities for $R \times C$ Contingency Tables [G2]

David L. March [Recd. 24 Nov. 1970 and 7 Mar. 1971]
School of Education, Lehigh University,
Bethlehem, PA. 18015

Key Words and Phrases: probability, contingency table, test of significance

CR Categories: 3.5, 5.5

Language: Fortran

Description

Freeman and Halton [1] derive a general method for computing exact probabilities for contingency tables that result if a sample is subjected to k different and independent classifications. The following algorithm is limited to the case where $k = 2$.

If a sample of size N is subjected to two different and independent classifications, A and B , with R and C classes respectively, the probability P_z of obtaining the observed array of cell frequencies $X(x_{ij})$, under the conditions imposed by the arrays of marginal totals $A(r_i)$ and $B(c_j)$ is given by

$$P_z = \frac{\prod_{i=1}^R (r_i!) \prod_{j=1}^C (c_j!)}{N! \prod_{i=1}^R \prod_{j=1}^C (x_{ij}!)} \quad (1)$$

Expression (1) is exact and holds if (a) the parent population is infinite or the sampling is done with replacement of the sampled items, (b) the sampling is random, (c) the population is homogeneous, and (d) the marginal totals are considered fixed in repeated sampling.

To test the null hypothesis that A and B are independent against the indefinite two-sided alternative, the probability P_s of obtaining an array as probable as, or less probable than, the observed array is needed. P_s is found as follows: (a) the probability P_t of the observed array is computed; (b) the probabilities for all other possible arrays of cell frequencies, subject to the conditions imposed by the fixed marginal totals, are computed; and (c) P_s is then obtained by

Copyright © 1972, Association for Computing Machinery, Inc.
General permission to republish, but not for profit, an algorithm is granted, provided that reference is made to this publication, to its date of issue, and to the fact that reprinting privileges were granted by permission of the Association for Computing Machinery.

summing all of the probability values found in (b) that are less than, or equal to, the probability P_t .

Method. The method of the subroutine uses the fact that expression (1) can be rewritten as

$$P_z = Q_z / R_z$$

where

$$Q_z = \frac{\prod_{i=1}^R (r_i!) \prod_{j=1}^C (c_j!)}{N!}$$

which is constant for the given set of marginal totals (r_i) and (c_j) and

$$R_z = \prod_{i=1}^R \prod_{j=1}^C (x_{ij}!)$$

which varies depending on the array of cell frequencies (x_{ij}) . In order to avoid machine overflow and roundoff error, these computations are performed using logarithms.

The observed $R \times C$ contingency table is specified by the $NR \times NC$ matrix which is partitioned as follows:

x_{11}	x_{1C}	r_1
\vdots				\vdots
x_{R1}	x_{RC}	r_R
c_1	c_C	N

After computing the constant term $QXLOG$ and the probability of the given table PT , the subroutine assigns to each of the lower right $(R - 1) \times (C - 1)$ cells the minimum of its corresponding row and column totals which is the maximum possible number for the cell. These cells are then varied in all possible combinations with each cell varied between its maximum number and zero.

Starting with cell $(2,2)$, the variation is accomplished by subtraction of 1. When the subtraction yields a zero or positive result the routine goes to compute the remainder of the cell frequencies. When a negative result is obtained, the cell in question, say cell (i,j) , is reset to the minimum of the corresponding row and column totals, 1 is subtracted from cell $(i, j + 1)$ or, if $j + 1$ is greater than C , cell $(i + 1, 2)$, and the count down resumes at cell $(2,2)$. If none of the lower right $(R - 1) \times (C - 1)$ cells yield a zero or positive result, the computations are complete and the subroutine returns to the caller. For example, if the top line (below) is the cell maximum ordered left to right from the $(2,2)$ to the (R, C) cell, the combinations generated will be

2	1	1	...
1	1	1	...
0	1	1	...
2	0	1	...
1	0	1	...
0	0	1	...
2	1	0	...
	\vdots		
0	0	0	...

The column 1 and row 1 cells are filled by subtraction of the generated cell numbers from the marginal totals. Since the method described above yields illegal as well as legal partitions, it is possible to obtain a negative result for one of these cells. When this occurs, the routine goes back to get a new set of cell frequencies. Otherwise

RXLOG is computed. Then, the probability *PX* is computed and added to the cumulative sum *PC*. If *PX* is less than, equal to, or, to avoid missing one due to computational inaccuracy, slightly larger than *PT*, *PX* is also added to the significance probability *PS*.

Since *PC* is the probability of obtaining some of the tables possible within the constraints of the marginal totals, *PC* should equal 1.0. The accuracy of the result can be estimated from the amount of deviation of *PC* from 1.0.

The floating point logarithms (base 10) of the integer factorials are obtained from function *FACLOG*. For arguments less than or equal to 100, the result is obtained from a table that is computationally filled on the first reference to *FACLOG*. Stirling's approximation is used for arguments greater than 100.

Results. The algorithm was tested on a CDC 6400 (60 bit word) using 2×3 ($N = 30$), 2×4 ($N = 7$), and 3×3 ($N = 7$) contingency tables. Results for the 2×3 tables were verified against values separately computed using programs developed by March [2]. In several cases *PC* deviated from 1.0. by 1.0×10^{-12} . Results for the 2×4 and 3×3 tests were verified by hand computation.

The author is indebted to the referees for their valuable comments and suggestions.

References

- Freeman, G.H., and Halton, J.H. Note on an exact treatment of contingency, goodness of fit, and other problems of significance. *Biometrika* 38 (1951), 141-149.
- March, D.L. Accuracy of the chi-square approximation for 2×3 contingency tables with small expectations. An unpublished D.Ed. Diss., School of Education, Lehigh U., Bethlehem, Pa., 1970.

Algorithm

```

SUBROUTINE CONP(MATRIX, NR, NC, PT, PS, PC)
C INPUT ARGUMENTS.
C MATRIX = SPECIFICATION OF THE CONTINGENCY TABLE.
C THIS MATRIX IS PARTITIONED AS FOLLOWS
C
C      X(11).....X(1C)  R(1)
C      .      .      .
C      .      .      .
C      X(R1).....X(RC)  R(R)
C      .      .      .
C      C(1).....C(C)   N
C
C WHERE X(IJ) ARE THE OBSERVED CELL FREQUENCIES,
C R(I) ARE THE ROW TOTALS, C(J) ARE THE COLUMN
C TOTALS, AND N IS THE TOTAL SAMPLE SIZE.
C NOTE THAT THE ORIGINAL CELL FREQUENCIES ARE
C DESTROYED BY THIS SUBROUTINE.
C
C NR = THE NUMBER OF ROWS IN MATRIX (R=NR-1).
C
C NC = THE NUMBER OF COLUMNS IN MATRIX (C=NC-1).
C
C OUTPUT ARGUMENTS.
C
C PT = THE PROBABILITY OF OBTAINING THE GIVEN TABLE.
C
C PS = THE PROBABILITY OF OBTAINING A TABLE AS PROBABLE
C AS, OR LESS PROBABLE THAN, THE GIVEN TABLE.
C
C PC = THE PROBABILITY OF OBTAINING SOME OF THE
C TABLES POSSIBLE WITHIN THE CONSTRAINTS OF THE
C MARGINAL TOTALS. (THIS SHOULD BE 1.0. DEVIATIONS
C FROM 1.0 REFLECT THE ACCURACY OF THE COMPUTATION.)
C
C EXTERNALS.
C
C FACLOG(N) = FUNCTION TO RETURN THE FLOATING POINT
C VALUE OF LOG BASE 10 OF N FACTORIAL.
C
C DIMENSION MATRIX(NR,NC)
C INTEGER R,C,TEMP
C
C R=NR-1
C C=NC-1
C
C COMPUTE LOG OF CONSTANT NUMERATOR
C
C QXLOG=-FACLOG(MATRIX(NR,NC))
C DO 10 I=1,R
C 10 QXLOG=QXLOG+FACLOG(MATRIX(I,NC))
C DO 20 J=1,C
C 20 QXLOG=QXLOG+FACLOG(MATRIX(NR,J))
C
C COMPUTE PROBABILITY OF GIVEN TABLE
C
C RXLOG=0.0
C DO 50 I=1,R
C DO 50 J=1,C
C 50 RXLOG=RXLOG+FACLOG(MATRIX(I,J))
C PT=10.0**(-QXLOG-RXLOG)

```

```

C
C PS=0.0
C PC=0.0
C
C FILL LOWER RIGHT (R-1) X (C-1) CELLS WITH
C MINIMUM OF ROW AND COLUMN TOTALS
C
C DO 100 I=2,R
C DO 100 J=2,C
C 100 MATRIX(I,J)=MINO(MATRIX(I,NC),MATRIX(NR,J))
C GO TO 300
C
C OBTAIN A NEW SET OF FREQUENCIES IN
C LOWER RIGHT (R-1) X (C-1) CELLS
C
C 200 DO 220 I=2,R
C DO 220 J=2,C
C MATRIX(I,J)=MATRIX(I,J)-1
C IF(MATRIX(I,J).GE.0) GO TO 300
C 220 MATRIX(I,J)=MINO(MATRIX(I,NC),MATRIX(NR,J))
C RETURN
C
C FILL REMAINDER OF OBSERVED CELLS
C .....COMPLETE COLUMN 1
C
C 300 DO 320 I=2,R
C TEMP=MATRIX(I,NC)
C DO 310 J=2,C
C 310 TEMP=TEMP-MATRIX(I,J)
C IF(TEMP.LT.0) GO TO 200
C 320 MATRIX(I,1)=TEMP
C
C .....COMPLETE ROW 1
C
C DO 340 J=1,C
C TEMP=MATRIX(NR,J)
C DO 330 I=2,R
C 330 TEMP=TEMP-MATRIX(I,J)
C IF(TEMP.LT.0) GO TO 200
C 340 MATRIX(1,J)=TEMP
C
C COMPUTE LOG OF THE DENOMINATOR
C
C RXLOG=0.0
C DO 350 I=1,R
C DO 350 J=1,C
C 350 RXLOG=RXLOG+FACLOG(MATRIX(I,J))
C
C COMPUTE PX. ADD TO PS IF PX .LE. PT
C (ALLOW FOR ROUND-OFF ERROR)
C
C PX=10.0**(-QXLOG-RXLOG)
C PC=PC+PX
C IF((PT/PX).GT.0.99999) PS=PS+PX
C GO TO 200
C END
C FUNCTION FACLOG(N)
C
C INPUT ARGUMENT.
C
C N = AN INTEGER GREATER THAN OR EQUAL TO ZERO.
C
C FUNCTION RESULT.
C
C FACLOG = THE LOG TO THE BASE 10 OF N FACTORIAL.
C
C DIMENSION TABLE(101)
C DATA TP106/0.39908 99342/
C DATA EL06 /0.43429 44819/
C DATA IFLAG/0/
C
C USE STIRLINGS APPROXIMATION IF N GT 100
C
C IF(N.GT.100) GO TO 50
C
C LOOK UP ANSWER IF TABLE WAS GENERATED
C
C IF(IFLAG.EQ.0) GO TO 100
C 10 FACLOG=TABLE(N+1)
C RETURN
C
C HERE FOR STIRLINGS APPROXIMATION
C
C 50 X=FLOAT(N)
C FACLOG=(X+0.5)*ALOG10(X) - X*EL06 + TP106
C 1 + EL06/(12.0*X) - EL06/(360.0*X*X*X)
C RETURN
C
C HERE TO GENERATE LOG FACTORIAL TABLE
C
C 100 TABLE(1)=0.0
C DO 120 I=2,101
C X=FLOAT(I-1)
C 120 TABLE(I)=TABLE(I-1)+ALOG10(X)
C IFLAG=1
C GO TO 10
C END

```