

Chapter g11 – Contingency Table Analysis

1. Scope of the Chapter

This chapter contains a function for the analysis of two-way contingency tables. Functions in Chapter g02 may be used to fit generalised linear models to discrete data.

2. Background

2.1. Discrete Data

Discrete data can be usefully categorized into three types.

- (1) *Binary data.* The variables can take one of two values, for example yes or no. The data may be grouped, for example the number of yes responses in ten questions.
- (2) *Categorical data.* The variables can take one of two or more values or levels, but the values are not considered to have any ordering; for example the values may be red, green, blue or brown.
- (3) *Ordered categorical data.* This is similar to categorical data but an ordering can be placed on the levels, for example: poor, average or good.

A second important categorization to be made is whether one of the discrete variables can be considered as a response variable or whether it is just the association between the discrete variables that is being considered. If the response variable is binary then a logistic or probit regression model can be used. These are special cases of the generalised linear model with binomial errors. Handling a categorical or ordered categorical response variable is more complex; for discussion of appropriate models see McCullagh and Nelder (1989).

To investigate the association between discrete variables a contingency table can be used.

2.2. Contingency Tables

The simplest case is the two-way table formed when considering two discrete variables. For a data set of n observations classified by the two variables with r and c levels respectively, a two-way table of frequencies or counts with r rows and c columns can be computed.

n_{11}	n_{12}	\dots	n_{1c}	$n_{1.}$
n_{21}	n_{22}	\dots	n_{2c}	$n_{2.}$
\vdots	\vdots	\vdots	\vdots	\vdots
n_{r1}	n_{r2}	\dots	n_{rc}	$n_{r.}$
$n_{.1}$	$n_{.2}$	\dots	$n_{.c}$	n

If p_{ij} is the probability of an observation in cell ij then the the model which assumes no association between the two variables is the model

$$p_{ij} = p_{i.}p_{.j}$$

where $p_{i.}$ is the marginal probability for the row variable and $p_{.j}$ is the marginal probability for the column variable, the marginal probability being the probability of observing a particular value of the variable ignoring all other variables. The appropriateness of this model can be assessed by two commonly used statistics:

the Pearson χ^2 -statistic

$$\sum_{i=1}^r \sum_{j=1}^c \frac{(n_{ij} - f_{ij})^2}{f_{ij}},$$

and the likelihood ratio test statistic

$$2 \sum_{i=1}^r \sum_{j=1}^c n_{ij} \times \log(n_{ij}/f_{ij}).$$

The f_{ij} are the fitted values from the model; these values are the expected cell frequencies and are given by

$$f_{ij} = n\hat{p}_{ij} = n\hat{p}_{i.}\hat{p}_{.j} = n(n_{i.}/n)(n_{.j}/n) = n_{i.}n_{.j}/n.$$

Under the hypothesis of no association between the two classification variables, both these statistics have, approximately, a χ^2 -distribution with $(c-1)(r-1)$ degrees of freedom. This distribution is arrived at under the assumption that the expected cell frequencies, f_{ij} , are not too small.

In the case of the 2 by 2 table, i.e., $c = 2$ and $r = 2$, the χ^2 -approximation can be improved by using Yates's continuity correction factor. This decreases the absolute value of $(n_{ij} - f_{ij})$ by $\frac{1}{2}$. For 2 by 2 tables with a small values of n the exact probabilities can be computed; this is known as Fisher's exact test.

An alternative approach, which can easily be generalised to more than two variables, is to use log-linear models. A log-linear model for two variables can be written as

$$\log(p_{ij}) = \log(p_{i.}) + \log(p_{.j}).$$

A model like this can be fitted as a generalised linear model with Poisson error with the cell counts, n_{ij} , as the response variable.

3. References

- Everitt B S (1977) *The Analysis of Contingency Tables* Chapman and Hall.
 Kendall M G and Stuart A (1969) *The Advanced Theory of Statistics (Volume 1)* Griffin (3rd Edition).
 McCullagh P and Nelder J A (1983) *Generalized Linear Models* Chapman and Hall.

4. Available Functions

- | | |
|---|--------|
| Computes χ^2 -statistics for a two-way contingency table | g11aac |
| The following routines may also be used to analyse discrete data: | |
| Fit a log-linear model to a contingency table | g02gcc |
| Fit generalized linear models to binary data | g02gbc |