

THE ANALYSIS OF CONTINGENCY TABLES

Log-linear Poisson models for weighted numbers

R-76-31

J. de Leeuw, Leyden State University, Dept. Data theory

S. Oppe, Institute for Road Safety Research SWOV

Voorburg, 1976

Institute for Road Safety Research SWOV, The Netherlands

CONTENTS

1. Introduction

2. Model
 - 2.1. Fundamental assumption
 - 2.2. Independence assumptions regarding characteristics in models
 - 2.3. Saturated and unsaturated models
 - 2.4. Weighted Poisson models

3. The design matrix
 - 3.1. General
 - 3.2. Three useful forms of design matrices
 - 3.2.1. Helmert matrices
 - 3.2.2. Orthogonal polynomials
 - 3.2.3. Between-within contrasts
 - 3.2.4. Combination of design matrices

4. Parameter estimation and hypothesis testing
 - 4.1. Introduction
 - 4.2. Modified minimum Chi-squared methods
 - 4.3. Calculations and limit distributions

Literature

- Annex 1: Correction for bias
Annex 2: Computer programme
Annex 3: Example of an analysis

1. INTRODUCTION

Contingency tables (or cross tables) classify elements of populations or samples (of varying kinds) with reference to one or more characteristics. For instance, classification of fatalities for a year according to age and mode of road usage. If there is only one characteristic, one often speaks of marginal tables. But this is also said of tables originating if one or more variables of a contingency table are added and one or more others are not. Since there is no essential difference between a marginal table and a contingency table (merely a functional difference), our future references will be to contingency tables.

The term "contingency table" is better than "cross table", because it expresses something of the assumptions made in cross-table analysis as regards the contingent factors assumed to play a part in creating the table. This aspect is essential in sampling.

By means of a sample we try on the one hand to describe the population from which the sample is taken, and on the other to verify opinions about this population. The models of analysis described below assume that a sample gives a more or less correct picture of the population, dependent only on random fluctuations.

Assumptions regarding the way chance plays a part form the basis of the model of analysis. Within it, different model specifications are again possible.

Analysis of contingency tables does not usually assume that there are specific relations (such as order relations or even metric relations) between the classes of a characteristic.

Such extra assumptions are possible, however, within specific models, for instance for a variable such as age.

In recent years there has been a new development in the way in which contingency tables are analysed. While it used to be customary (mostly by the Chi-squared test) to verify overall-hypotheses about a table with either one or two characteristics, analysis now increasingly stresses the detailed information the table contains. Furthermore, it is also possible to analyse higher-order tables (subdivided into a number of characteristics), in order that more complex relationships, i.e. relations between more than two characteristics at once, can be investigated.

2. MODEL

2.1. Fundamental assumption

The fundamental assumption is that the number of accidents in the cells of the contingency table are independent random variables with a Poisson distribution, in which the Poisson distribution parameters may differ. To keep this fairly concrete: if we are concerned with a two-way table with r rows and k columns, then we could write the Poisson assumption for each cell as follows:

there are numbers $\lambda_{ij} \geq 0$ ($i = 1, \dots, r; j = 1, \dots, k$) such that:

$$\text{prob} \left[\begin{matrix} X_{ij} \\ \sim_{ij} \end{matrix} = x_{ij} \right] = e^{-\lambda_{ij}} \frac{\lambda_{ij}^{x_{ij}}}{x_{ij}!}$$

In this, X_{ij} is the stochastic variable of cell (i,j) which can assume as values the natural numbers $x_{ij} = 0, 1, 2, \dots$. A briefer way of writing this assumption is:

$$X_{ij} \sim \mathcal{P}(\lambda_{ij}),$$

which we can read as: X_{ij} has a Poisson distribution with a parameter λ_{ij} .

2.2. Independence assumptions regarding characteristics in models

Although we assume that the X_{ij} are independent, it is of course possible that there exist relations between parameters λ_{ij} . By investigating the relations between these parameters we can examine whether the characteristics the variables possess are also independent. What do we mean when we say that the rows and columns of an $r \times k$ contingency table (with independent Poisson variables X_{ij}) correspond to independent row and column variables? Suppose $X_{i.}$ and $X_{.j}$ are the marginal distributions, viz:

$$X_{i.} = \sum_{j=1}^k X_{ij}$$

and

$$X_{.j} = \sum_{i=1}^r X_{ij}$$

The requirement that the row and column variables must be independent means that the chances of the r conditional distributions within rows:

$$\text{prob} \left[(X_{i1} = x_{i1}) \wedge (X_{i2} = x_{i2}) \wedge \dots \wedge (X_{ik} = x_{ik}) \mid X_{i.} = x_{i.} \right]$$

are the same for all $i=1, \dots, r$, and that the k conditional distributions within columns:

$$\text{prob} \left[(X_{1j} = x_{1j}) \wedge (X_{2j} = x_{2j}) \wedge \dots \wedge (X_{rj} = x_{rj}) \mid X_{.j} = x_{.j} \right]$$

are the same for all $j=1, \dots, k$. Using the independence of the X_{ij} and the Poisson assumption, we can infer that the conditional distributions within rows are the same as the multinomial distributions:

$$\frac{x_{i.}!}{\prod_{j=1}^k x_{ij}!} \cdot \prod_{j=1}^k \left(\frac{\lambda_{ij}}{\lambda_{i.}} \right)^{x_{ij}},$$

while the conditional distributions within columns are the same as the multinomial distributions:

$$\frac{x_{.j}!}{\prod_{i=1}^r x_{ij}!} \cdot \prod_{i=1}^r \left(\frac{\lambda_{ij}}{\lambda_{.j}} \right)^{x_{ij}}$$

The row and column variables are thus independent when $\left(\frac{\lambda_{ij}}{\lambda_{i.}}\right)$ is the same for all j , and $\left(\frac{\lambda_{ij}}{\lambda_{.j}}\right)$ is the same for all i .

Necessary and sufficient conditions for this are that there are numbers $a_i \geq 0$ ($i=1, \dots, r$), $b_j \geq 0$ ($j=1, \dots, k$), and $m \geq 0$, such that $\lambda_{ij} = m a_i b_j$ for all i, j .

This multiplicative model is mostly converted into a linear model by taking the logarithm:

$$\ln \lambda_{ij} = \mu + \alpha_i + \beta_j \quad (1)$$

in which $\alpha_i = \ln a_i$ and so on.

Hence, such models are also called log-linear models. The log-linear model (1) is therefore equivalent with the requirement of independence of the row and column variables.

2.3. Saturated and unsaturated models

As stated, in addition to testing hypotheses regarding tables, description of the tables is sometimes also required. If the characteristics are not independent and the above model (1) does not therefore apply, it can be extended with specific parameters for the cells. In that case we have the following model

$$\ln \lambda_{ij} = \mu + \alpha_i + \beta_j + \gamma_{ij} \quad (2)$$

With regard to this model it is always possible to find such parameters μ , α_i , β_j and γ_{ij} that there is complete agreement between the table one wishes to describe and the model used for description. The significance of the description is that the variation in the numbers of observations of the cells in the table is shown in relation to the structure of the table: it can be seen, for instance, to what extent the variation results from a row-effect, a column-effect or an interaction effect. Although there are now as many parameters as cells and

hence there is no reduction in information, there is an ordering of information. Moreover, note that model (1) is a special case of model (2); it is the same, except for the restriction that $\gamma_{ij} = 0$ for all i, j . Other restrictions are also possible, for instance that the α_i 's together form a linear relation or, for example, are equal to zero. In all such cases we speak of unsaturated models. If we are dealing with a sample, these non-saturated models can be regarded as verifiable hypotheses regarding the population from which the sample originates. With a saturated model, such verification is not possible because the model fully describes the data.

As regards the choice of the model of analysis, there is close agreement with linear models as used in analysis of variance. Here, again, we can speak of a breakdown of the table into components: how great is the row-contribution, the column-contribution, the unique cell contribution of each cell? For an incidental table this can be examined by estimating the parameters of the model.

This systematic breakdown therefore provides an efficient review of the information contained in the table. It is also possible to give confidence limits of the estimators for the parameters, so that verification of individual estimators is also possible.

A good description of the relation between analysis of variance models and log-linear models is given by Nelder & Wedderburn (1972).

2.4. Weighted Poisson models

So far we treated only the numbers of accidents as a function of a number of characteristics. But we are sometimes interested in analysing accident figures normalised for a given exposure factor such as number of inhabitants, road lengths, and so on. If we enter the numbers of accidents in the table with a measure of exposure per cell, which may differ from cell to cell, we can use a more general Poisson model. The fundamental assumption now becomes

$$x_{ij} \sim P(c_{ij} \lambda_{ij}),$$

in which the c_{ij} are the given exposure factors, and in which a log-linear model is again assumed for the λ_{ij} .

3. THE DESIGN MATRIX

3.1. General

In matrix notation, the general form of a log-linear model for n Poisson variables $\mathbf{X}_1 \sim P(\lambda_1)$ can be written as

$$\eta = V\theta,$$

in which η is a vector of values $\eta_1 = \ln \lambda_1$. V is a given matrix of the order $n \times p$ (known as the design matrix), and θ a vector of p unknown parameters. If the \mathbf{X}_1 are arranged in a two-way table and if we replace the index 1 by the row and column indices i and j , then we can rewrite the model

$$\eta_{ij} = \ln \lambda_{ij} = \mu + \alpha_i + \beta_j$$

when $r = k = 2$, for instance as

$$\begin{bmatrix} \eta_{11} \\ \eta_{21} \\ \eta_{12} \\ \eta_{22} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

Note that in this case the design matrix V is of the order 4×5 and rank 3. This become clear if we rewrite the model in the equivalent form:

$$\begin{bmatrix} \eta_{11} \\ \eta_{21} \\ \eta_{12} \\ \eta_{22} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & -1 \\ 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

with:

$$\begin{aligned}\theta_1 &= \mu + \bar{\alpha} + \bar{\beta} \\ \theta_2 &= (\alpha_1 - \bar{\alpha}) = -(\alpha_2 - \bar{\alpha}) \\ \theta_3 &= (\beta_1 - \bar{\beta}) = -(\beta_2 - \bar{\beta})\end{aligned}$$

in which $\bar{\alpha}$ and $\bar{\beta}$ respectively are the means of the α 's and β 's. Generally, it is always possible (and advisable) to choose the design matrix so that its rank is the same as its number of columns. This obviates having to impose extra restrictions on the parameters in order to find a unique solution. If we were to seek a direct solution for the α 's and β 's, these restrictions would be:

$$\alpha_1 + \alpha_2 = 0 \text{ and } \beta_1 + \beta_2 = 0.$$

The rank matrix of the columns and the number of columns are, for instance, always the same if we choose V so that $V'V$ is diagonal, with V' as the transpose of matrix V (V is then called column-wise orthogonal), or so that $V'V$ is equal to the unit matrix (V is then called column-wise orthonormal).

3.2. Three useful forms of design matrices

3.2.1. Helmert matrices

Let us first consider the case in which we have a single classification. Example: $i=1, \dots, n$ corresponds to n age categories, $X_{\sim i}$ is the number of accidents in each such category. A first-type design matrix often used is the Helmert matrix.

A complete Helmert matrix for $n = 4$ is as follows:

$$\begin{array}{cccc} 1 & -1 & -1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & 0 & 2 & -1 \\ 1 & 0 & 0 & 3 \end{array}$$

Note that this V is column-wise orthogonal. The model $\eta = V\theta$ is therefore saturated. A perfect fit is possible if we choose

$\hat{\theta} = (V'V)^{-1}V'r$. Unsaturated models are possible by omitting columns of V , which agrees with the hypothesis that some of the elements of θ in the saturated model are equal to zero. The interpretation of Helmert effects becomes clear from the following equivalences:

$$\theta_1 = 0 \Leftrightarrow \sum r_i = 0 \Leftrightarrow \sqrt[n]{\prod_{i=1}^n \lambda_i} = 1$$

$$\theta_2 = 0 \Leftrightarrow r_2 = r_1 \Leftrightarrow \lambda_2 = \sqrt{\lambda_1}$$

$$\theta_3 = 0 \Leftrightarrow 2r_3 = r_1 + r_2 \Leftrightarrow \lambda_3 = \sqrt{2\lambda_1\lambda_2}$$

$$\theta_4 = 0 \Leftrightarrow 3r_4 = r_1 + r_2 + r_3 \Leftrightarrow \lambda_4 = \sqrt[3]{\lambda_1\lambda_2\lambda_3}$$

From this, we can for instance derive:

$$\theta_3 = \theta_4 = 0 \Leftrightarrow \lambda_3 = \lambda_4 = \sqrt{2\lambda_1\lambda_2},$$

and so on. Helmert effects therefore compare every λ_i individually with the geometric mean of the preceding λ_i . In this way, we can discover whether there is a trend in our data, or perhaps a sudden jump.

3.2.2. Orthogonal polynomials

Let us assume that the age categories in our example are intervals of equal length. We might then be interested in the functional relation between age and accident rate. We can describe this functional relation as a polynomial, i.e. as a linear combination of orthogonal polynomials; for $n = 3$ this gives, for instance, the following (column-wise orthogonal) design matrix:

$$\begin{array}{ccc} 1 & -1 & 1 \\ 1 & 0 & -2 \\ 1 & 1 & 1 \end{array}$$

Each constant function on (1,2,3) is of course a multiple of the first column of V, each linear function on (1,2,3) is a linear combination of the first two columns, and each second-degree function is a linear combination of the first three columns. Each function on (1,2,3) can be regarded as a second degree function: this is merely another way of saying that the model defined by V is saturated. Unsaturated models generally have the form $\theta_3 = 0$ or $\theta_2 = \theta_3 = 0$. The hypothesis $\theta_3 = 0$ says that the three points $(1, \eta_1)$, $(2, \eta_2)$, and $(3, \eta_3)$ are on a straight line, the hypothesis $\theta_2 = \theta_3 = 0$ says that $\eta_1 = \eta_2 = \eta_3$. Generally, the hypothesis that (η_1, \dots, η_n) is a q^{th} degree polynomial on $(1, 2, \dots, n)$ can be written $\eta_i = \prod_q(i)$. From our discussion it follows that:

$$\eta_i = \prod_q(i) \Leftrightarrow \theta_{q+1} = \dots = \theta_n = 0.$$

The interpretation of polynomial effects in log-linear models is rendered difficult by the use of log-transformation, since:

$$\eta_i = \prod_q(i) \Leftrightarrow \lambda_i = \exp(\prod_q(i)) = \exp(\alpha_0 + \alpha_1 i + \dots + \alpha_q i^q) = \left[\exp(i^0) \right]^{\alpha_0} \left[\exp(i^1) \right]^{\alpha_1} \dots \left[\exp(i^q) \right]^{\alpha_q}.$$

This latter function is rather less simple and unacquainted than a polynomial.

3.2.3. Between-Within contrasts

In many cases, categories of our classification break down naturally into different groups. Age can, for instance, be divided into two groups: the under-forties and over-forties. This subdivision can be shown in saturated design-matrix form as

00-20	1	-1	-1	0
20-40	1	-1	+1	0
40-60	1	1	0	-1
60-80	1	1	0	+1

In this case the measurements themselves are therefore in four categories, and we examine so to speak whether subdivision into fewer categories is possible without forfeiting too much information. The first column of V corresponds as usual to the total average, the second column contrasts the two groups (the between-group effect), and the third and fourth columns examine the effects within the groups individually. If there are K groups with n_k elements ($\sum_{k=1}^K n_k = n$), then there are generally K-1 between-group effects, and $\sum_{k=1}^K (n_k - 1) = n - K$ within-group effects. The most common unsaturated models state that all θ values corresponding to between-group effects are zero. This agrees with the hypothesis that the arithmetic means of the η_i are the same for each group, which is equivalent to the fact that the geometric means of the λ_i are the same for every group.

3.2.4. Combination of design matrices

Let us now examine a two-way classification with, for instance, two classes in the first characteristic (e.g. male against female), and four classes in the second characteristic (e.g. the four age groups in the preceding section). We first choose two design matrices V_1 and V_2 for the separate characteristics.

For example:

$$V_1 = \begin{matrix} +1 & -1 \\ +1 & +1 \end{matrix}$$

$$V_2 = \begin{matrix} +1 & -1 & -1 & 0 \\ +1 & -1 & +1 & 0 \\ +1 & +1 & 0 & -1 \\ +1 & +1 & 0 & +1 \end{matrix}$$

We next form from all $2 \times 4 = 8$ combinations of columns of V_1 and V_2 the external product (the external product of an n-vector x and an m-vector y is an n x m matrix with the elements $x_i y_j$). This gives the following eight matrices:

V_1	V_2	Product			
1	1	+1	+1	+1	+1
		+1	+1	+1	+1
1	2	-1	-1	+1	+1
		-1	-1	+1	+1
1	3	-1	+1	0	0
		-1	+1	0	0
1	4	0	0	-1	+1
		0	0	-1	+1
2	1	-1	-1	-1	-1
		+1	+1	+1	+1
2	2	+1	+1	-1	-1
		-1	-1	+1	+1
2	3	+1	-1	0	0
		-1	+1	0	0
2	4	0	0	+1	-1
		0	0	-1	+1

We can treat these eight matrices as eight vectors of eight elements, and thus form a design matrix V_{12} with these vectors as columns.

Thus:

Design matrix

+1	-1	-1	0	-1	+1	+1	0
+1	-1	+1	0	-1	+1	-1	0
+1	+1	0	-1	-1	-1	0	+1
+1	+1	0	+1	-1	-1	0	-1
+1	-1	-1	0	+1	-1	-1	0
+1	-1	+1	0	+1	-1	+1	0
+1	+1	0	-1	+1	+1	0	-1
+1	+1	0	+1	+1	+1	0	+1

Belonging to vector

n_{11}
n_{12}
n_{13}
n_{14}
n_{21}
n_{22}
n_{23}
n_{24}

The matrix V_{12} so constructed is again column-wise orthogonal, and defines a saturated model. We can say that V_{12} is formed via external products. In using a design matrix built up in this way, we usually wish to investigate a given type of unsaturated models. Let us examine these unsaturated models with respect to our example.

We first choose the column corresponding to the first column of V_1 and the first column of V_2 . This is the first column of V_{12} . The hypothesis $\theta_1 = 0$ is equivalent to the hypothesis that the arithmetic mean of the η_{ij} ($i=1,2;j=1,2,3,4$) is zero, i.e. that the geometric mean of the λ_{ij} equals one.

We next choose the group of columns of V_{12} composed from the first column of V_1 and column two, three or four of V_2 . These are columns 2, 3, 4 of V_{12} . The hypothesis $\theta_2 = \theta_3 = \theta_4 = 0$ is equivalent to the hypothesis that the column averages of the η_{ij} are identical, or:

$$\eta_{.1} = \eta_{.2} = \eta_{.3} = \eta_{.4}$$

This is equivalent to:

$$\lambda_{11} \lambda_{12} = \lambda_{12} \lambda_{22} = \lambda_{13} \lambda_{23} = \lambda_{14} \lambda_{24}$$

In the same way, we can choose the group of columns composed from the first column of V_2 and a non-first column of V_1 . This group consists of the fifth column of V_{12} . The hypothesis $\theta_5 = 0$ is:

$$\lambda_{11} \lambda_{12} \lambda_{13} \lambda_{14} = \lambda_{21} \lambda_{22} \lambda_{23} \lambda_{24}$$

Lastly, there is the group of columns 6, 7, 8 corresponding to a non-first column of V_1 and a non-first column of V_2 . The hypothesis $\theta_6 = \theta_7 = \theta_8 = 0$ corresponds to:

$$\eta_{ij} = \frac{\eta_{i.}}{k} + \frac{\eta_{.j}}{r} - \frac{\eta_{..}}{kxr}$$

that is to say with lack of additive interaction in the η_{ij} (cf. model (1) on page 6), which in turn is the same as the lack of multiplicative interaction in the λ_{ij} (for a comparison of these

two forms of interaction see Darroch, 1974; Lancaster, 1973, 1975). It is clear that this form of analysis via external products can be generalised to tables with more than two classifications. We always begin with design matrices for each of the characteristics, form external products, and group the columns of the ultimate design matrix by examining which first columns appear in them. Hence, we form groups of effects corresponding to the additive interactions of the η 's (which are known from ordinary analysis of variance, and to multiplicative interactions of the λ 's (which can be interpreted in the manner of section 2.2. as independence models). It is important to realise that an interaction hypothesis in the form $\theta_6 = \theta_7 = \theta_8 = 0$ in the above example is either true or not true, regardless of the choice of the original $V_1, V_2 \dots$. The choice of design matrix for a given characteristic is therefore of importance only for better interpretation of the individual θ 's, but is of no importance in describing the table according to the contributions of the characteristics or the interactions between them.

4. PARAMETER ESTIMATION AND HYPOTHESIS TESTING

4.1. Introduction

For convenience, we will briefly enumerate the fundamental assumptions for the class of models in which we are interested.

$$A1: \tilde{x}_i \sim P(\ell_i, \lambda_i^o)$$

A2: \tilde{x}_i are independent

$$A3: \eta^o = V \theta^o.$$

In A1, ℓ is a known vector of weights (or measures of exposure); in A3, $\eta^o = \ln \lambda^o$, and V is a known $n \times p$ design matrix, which we shall assume to be column-wise orthonormal. The superscript 'o' with θ , η , and λ is to indicate the 'actual' value of these parameters, and to distinguish them from estimators and variables in specific functions. What interests us in the first place is estimation of the p unknown parameters, and in the second place verifying whether the model A1, A2, A3 is correct. For this, it is important also to formulate in A3 other (equivalent) ways. If V is an $n \times p$ column-wise orthonormal matrix, then there is an $n \times (n - p)$ column-wise orthonormal matrix V_c such that $V'V_c = 0$. It is clear that A3 can also be written as:

$$A3: V_c' \eta^o = 0.$$

A third formulation is possible if we define the p -dimensional linear space \mathcal{U} as:

$$\mathcal{U} = \{ \eta \mid V_c' \eta = 0 \} = \{ \eta \mid \eta = V \theta \}.$$

then

$$A3: \eta^o \in \mathcal{U}.$$

It is generally unfeasible to use estimators and test procedures which are optimal for all conceivable sample sizes. We shall therefore use asymptotic arguments, and derive estimates and tests which have optimum properties if certain factors tend to infinity. For this purpose, we reformulate A1 as:

$$A1: \underset{\sim i}{X}_i \sim P(m c_i \lambda_i^0).$$

The factor m indicates how great our weights c_i and parameters λ_i^0 are on average. If we continue our observations, the $\underset{\sim i}{X}_i$ will of course tend to infinity. The assumption A1 says, in fact, that all $\underset{\sim i}{X}_i$ tend to infinity just as quickly: if m becomes infinitely great, then the values $\underset{\sim i}{X}_i/m$ converge (in probability) towards the fixed factors $c_i \lambda_i^0$.

For our analyses, it is generally unnecessary to know the value of m ; we must merely be prepared to make this assumption. The following facts are known from the general theory of asymptotic statistical analysis. In the first place we shall be interested in estimators that are consistent; viz. if $m \rightarrow \infty$ then $\hat{\theta}(m) \xrightarrow{P} \theta^0$. In the second place, we are interested in estimators that are asymptotically normal, which means that their distribution more and more resembles a multi-normal distribution if m tends to infinity. For estimators with these two properties, which we can summarise as:

$$T1: m^{\frac{1}{2}} (\hat{\theta}(m) - \theta^0) \xrightarrow{L} \mathcal{N}(0, \Sigma)$$

the asymptotic dispersion matrix Σ satisfies the dissimilarity:

$$\Sigma \geq (V'M^0V)^{-1}$$

in which M^0 is the diagonal matrix with the values $c_i \lambda_i^0$ on the diagonal. Estimators in this class for which the dissimilarity is a similarity, and which in a certain sense are thus as precise as possible, are called efficient. Although nearly all available estimators satisfy T1, they do not necessarily meet the stricter requirement:

$$T2: m^{\frac{1}{2}} (\hat{\theta}(m) - \theta^0) \xrightarrow{L} \mathcal{N}(0, (V'M^0V)^{-1}).$$

Since efficiency is a necessary property, we shall confine ourselves to efficient estimators (i.e. estimators satisfying T2). Moreover, it is important to add that confidence intervals of estimators and tests of hypotheses via these estimators are generally asymptotically optimal if the estimators are efficient.

It is known that efficient estimators can be found by maximising the likelihood function, that gives the likelihood of the observations as a function of the parameters, and that an asymptotically optimum test of A3 within A1-A2 is possible by calculating the likelihood ratio between the most suitable estimator or estimators and the hypothetical value of the parameter or parameters. The estimation and test theory based on this maximum likelihood is set forth for log-linear Poisson models in Haberman (1974). The theory is modified for weighted Poisson models in De Leeuw (1975). On the whole, calculations based on likelihood are not very simple, and we consider here a different class of estimators and tests (also optimal and efficient), based on Neyman's modified minimum Chi-squared method (1949).

4.2. Modified minimum Chi-squared methods

We begin this section with a known limit theorem for Poisson variables which, applied to A1, says that if $m \rightarrow \infty$:

$$\frac{X_i - m e_i \lambda_i^0}{(m e_i \lambda_i^0)^{\frac{1}{2}}} \xrightarrow{L} \mathcal{N}(0, 1)$$

If we define

$$\tilde{Y}_i = \tilde{X}_i / (m e_i)$$

then we can rewrite it in a rather more convenient form:

$$m^{\frac{1}{2}} (\tilde{Y}_i - \lambda_i^0) \xrightarrow{L} \mathcal{N}(0, \frac{\lambda_i^0}{e_i}).$$

If, lastly, we define:

$$z_i = \ln \frac{Y_i}{\tilde{z}_i}$$

it follows from this that

$$m^{\frac{1}{2}} (z_i - \eta_i^0) \xrightarrow{L} \mathcal{N}(0, (e_i \lambda_i^0)^{-1})$$

The modified minimum Chi-squared method we shall discuss below has a simple geometric interpretation. We define the distance measure:

$$\delta(\eta_1, \eta_2) = (\eta_1 - \eta_2)' \tilde{X} (\eta_1 - \eta_2).$$

The matrix \tilde{X} is diagonal, and the X_i are on the diagonal. Note that so far we have already demonstrated that:

$$\delta(z, \eta^0) \xrightarrow{L} \chi_n^2$$

if $m \rightarrow \infty$ (this follows from the limit distribution of Z , and from $X_i/m \xrightarrow{P} e_i \lambda_i^0$). For estimations, we consider the distance between the vector Z of observations, and the collection of permitted estimators η . For calculating the modified minimum Chi-squared estimators we must choose $\hat{\eta}$ so that:

$$\delta(z, \hat{\eta}) = \min_{\eta \in \mathcal{V}} \delta(z, \eta).$$

This gives an estimation $\hat{\eta}$ for η^0 . The corresponding estimator for θ^0 is $V\hat{\eta}$, and the statistic used for testing A3 is $\delta(z, \hat{\eta})$. In the next section we study the distribution of estimators and statistics.

4.3. Calculations and limit distributions

The problem

$$\min_{\eta \in \mathcal{V}} \delta(z, \eta)$$

can be formulated in two different ways. The first formulation is:

$$\min_{\theta} \delta(z, v\theta)$$

This gives estimators $\hat{\theta}_I$, and then $\eta_I = V\hat{\theta}_I$. Formulation II uses Lagrange multipliers and can be written as

$$\min_{\eta} \max_{\omega} \delta(z, \eta) + 2\omega' v' \eta.$$

This gives estimators $\hat{\eta}_{II}$ and $\hat{\omega}$, and then $\hat{\theta}_{II} = v' \hat{\eta}_{II}$. The $(n-p)$ -element vector ω is a vector of indefinite multipliers. The model can now also be written as:

$$A3 : \omega = 0.$$

As the solution of the original problem is unique, obviously

$$\hat{\theta}_I = \hat{\theta}_{II} = \hat{\theta} \text{ and}$$

$$\hat{\eta}_I = \hat{\eta}_{II} = \hat{\eta}.$$

It follows from formulation I that $\hat{\theta}$ is given by

$$\hat{\theta} = (v'XV)^{-1}v'XZ$$

and hence

$$\hat{\eta} = v(v'XV)^{-1}v'XZ.$$

It further follows that both $\hat{\theta}$ and $\hat{\eta}$ are efficient estimators; in other words:

$$\begin{aligned} m^{\frac{1}{2}} (\hat{\theta} - \theta^0) &\xrightarrow{L} \mathcal{N}(0, (v'M^0v)^{-1}) \\ m^{\frac{1}{2}} (\hat{\eta} - \eta^0) &\xrightarrow{L} \mathcal{N}(0, v(v'M^0v)^{-1}v') \end{aligned}$$

The asymptotic dispersion matrices can be estimated by:

$$s(\hat{\theta}) = (V'XV)^{-1}$$

$$s(\hat{\eta}) = V(V'XV)^{-1}V'$$

It furthermore follows from the given results that:

$$\delta(z, \hat{\eta}) \xrightarrow{L} \chi^2_{n-p}.$$

Formulation II gives other useful information. We find:

$$\hat{\omega} = (V_c'X_c^{-1}V_c)^{-1}V_c'Z_c$$

$$\hat{\eta} = Z_c - X_c^{-1}V_c(V_c'X_c^{-1}V_c)^{-1}V_c'Z_c.$$

The vectors $\hat{\omega}$ and $\hat{\eta}$ are asymptotically independent, and

$$m^{\frac{1}{2}} \hat{\omega} \xrightarrow{L} \eta(0, (V_c'M_o^{-1}V_c)^{-1})$$

From equations I and II it also follows that $\delta(z, \hat{\eta})$ can be written in three different forms:

$$\begin{aligned} \delta(z, \hat{\eta}) &= Z_c' \left[X_c - X_cV(V'XV)^{-1}V'X_c \right] Z_c \\ &= \hat{\omega}' V_c' X_c^{-1} V_c \hat{\omega} \\ &= Z_c' V_c (V_c' X_c^{-1} V_c)^{-1} V_c' Z_c. \end{aligned}$$

The statistic $\delta(z, \hat{\eta})$ is therefore also found if we test A3 in the form $V_c'\eta = 0$ or $\omega = 0$ by using the asymptotic distribution of $V_c'Z_c$ and $\hat{\omega}$. These tests are known respectively as the Wald test and the Lagrange multiplier test; in this context they are thus equivalent to the Neyman method.

Especially if V_c is a matrix of low-rank, the Wald test will be preferable.

LITERATURE

1. J.N. Darroch. 'Multiplicative and additive interaction in contingency tables'. Biometrika, 1974, p. 207.
2. L.A. Goodman. 'The Multivariate Analysis of Qualitative Data: Interactions Among Multiple Classifications'. J.A.S.A., 1970, p. 226.
3. L.A. Goodman. 'Guided and Unguided Methods for the Selection of Models for a Set of T Multidimensional Contingency Tables'. J.A.S.A., 1973, p. 165.
4. S.J. Haberman. 'The Analysis of Frequency Data'. Univ. of Chicago Press, 1974.
5. H.O. Lancaster. 'The multiplicative definition of interaction'. Austral. J. Statist., 1971, p. 36.
6. J. de Leeuw. 'Maximum Likelihood Estimation for Weighted Poisson Models'. RN005-75. Leyden State University, Dept. Data theory, 1975.
7. J.A. Nelder & R.W.M. Wedderburn. 'Generalized Linear Models'. J.R. Statist. Soc. A, 1972, p. 370.
8. J. Neyman. 'Contributions to the theory of the X^2 -test'. Proc. of the Berkely Symp. on Math., Statist. and Probability, 1949, p. 239.
9. R.L. Plackett. 'The Analysis of Categorical Data'. Griffin, London, 1974.

ANNEX 1. CORRECTION FOR BIAS

We can apply a correction for bias if, before calculating the Z_i values, we first add $\frac{1}{2}$ to the X-values, so that Z_i is then defined as:

$$Z_i = \ln \frac{X_i + \frac{1}{2}}{m e_i}$$

Why $\frac{1}{2}$. Suppose we define:

$$\tilde{Z} = \ln \frac{\tilde{X} + a}{m \tilde{e}}$$

(For convenience we omit below the i and the small superscript o ; we also define $\mu = e^\lambda$). If:

$$\tilde{U} = \frac{(\tilde{X} - m \mu) + a}{m \mu}$$

then

$$\tilde{Z} = \eta + \tilde{U} - \frac{1}{2} \tilde{U}^2 + \frac{1}{3} \tilde{U}^3 - \frac{1}{4} \tilde{U}^4 + \dots$$

From which it follows that:

$$E(\tilde{Z}) = \eta + \frac{1}{m} \left(\frac{2a - 1}{\mu} \right) - \left(\frac{1}{m} \right)^2 \left(\frac{6a^2 - 12a + 5}{12 \mu^2} \right) + o(m^{-2}).$$

This correction also has the useful side-effect that Z is now also defined for $X = 0$.

ANNEX 2. COMPUTER PROGRAMME

In the computer programme it is necessary to read in per variable a design matrix of orthogonal column vectors, the first column vector being generated. The definitive design matrix is constructed in the programme with the aid of the external product method and then converted into an orthonormal matrix. If one is not interested in individual effects, therefore, the simplest method is to introduce Helmert effects. The θ 's of the saturated model are calculated with the formula:

$$\theta = (\underset{\sim}{V}'\underset{\sim}{X}\underset{\sim}{V})^{-1}\underset{\sim}{V}'\underset{\sim}{X}\underset{\sim}{Z}$$

In the case of a saturated model for the orthonormal V-matrix this formula reduces to:

$$\theta = \underset{\sim}{V}'\underset{\sim}{X}^{-1}\underset{\sim}{V}\underset{\sim}{V}'\underset{\sim}{X}\underset{\sim}{Z} = \underset{\sim}{V}'\underset{\sim}{Z}$$

The relevant variances, on the basis of which the standard scores are calculated, are on the diagonal of matrix $(\underset{\sim}{V}'\underset{\sim}{X}\underset{\sim}{V})^{-1}$, which is calculated in the case of saturation as $\underset{\sim}{V}'\underset{\sim}{X}^{-1}\underset{\sim}{V}$, and hence no inversion is needed.

For testing hypotheses in which (always limited) groups of θ 's are taken as zero. Formulation II on page 20 is used because in this case only a matrix of limited order need be inverted in order to obtain:

$$\delta (\underset{\sim}{z}, \hat{\underset{\sim}{\eta}}) = \underset{\sim}{z}'\underset{\sim}{V}_c (\underset{\sim}{V}'_c \underset{\sim}{X}^{-1}\underset{\sim}{V}_c)^{-1}\underset{\sim}{V}'_c \underset{\sim}{z}$$

The matrix $\underset{\sim}{V}'_c \underset{\sim}{X}^{-1}\underset{\sim}{V}_c$ is given as a part matrix of matrix $\underset{\sim}{V}'\underset{\sim}{X}^{-1}\underset{\sim}{V}$ already calculated.

ANNEX 3. EXAMPLE OF AN ANALYSIS

As an illustration, an example is worked out below.

It has been chosen because of the simplicity of the table.

A three-way table was chosen, in which the variables are:

A: The province of Noord-Brabant as against the Rest of the Netherlands.

B: Drinking established as against not established.

C: Location on road (intersection, road section, corner/bend).

The cells of the table show the number of deaths in the years 1971-1973 (Central Statistical Office data), inside built-up areas.

		C ₁ (intersection)	C ₂ (straight road)	C ₃ (corner/bend)
A ₁ (N-Br)	B ₁ (drinking)	22	48	14
	B ₂ (not drinking)	243	272	48
A ₂ (Rest of Neth.)	B ₁	97	202	68
	B ₂	1206	1442	189

These figures are weighted in the analysis for number of inhabitants in Noord-Brabant by factor 18.80 and in the Rest of the Netherlands by factor 115.08.

In analysis, use was made of the following design matrix V , built up from Helmert-effects.

Matrix:

$$V' = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & -1 \\ 1 & 1 & 1 & -1 & -1 & -1 & 1 & 1 & 1 & -1 & -1 & -1 \\ 1 & -1 & 0 & 1 & -1 & 0 & 1 & -1 & 0 & 1 & -1 & 0 \\ 1 & 1 & -2 & 1 & 1 & -2 & 1 & 1 & -2 & 1 & 1 & -2 \\ 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 \\ 1 & -1 & 0 & 1 & -1 & 0 & -1 & 1 & 0 & -1 & 1 & 0 \\ 1 & 1 & -2 & 1 & 1 & -2 & -1 & -1 & 2 & -1 & -1 & 2 \\ 1 & -1 & 0 & -1 & 1 & 0 & 1 & -1 & 0 & -1 & 1 & 0 \\ 1 & 1 & -2 & -1 & -1 & 2 & 1 & 1 & -2 & -1 & -1 & 2 \\ 1 & -1 & 0 & -1 & 1 & 0 & -1 & 1 & 0 & 1 & -1 & 0 \\ 1 & 1 & -2 & -1 & -1 & 2 & -1 & -1 & 2 & 1 & 1 & -2 \end{bmatrix}$$

Effect:

- T: total
- A: Noord-Brabant as against Rest of Netherlands
- B: drinking established as against not established
- C₁: intersection as against road section
- C₂: intersection + road section as against corner/bend
- A x B
- A x C₁
- A x C₂
- B x C₁
- B x C₂
- A x B x C₁
- A x B x C₂

The results of the estimations for the saturated model are given below. The (2x2x3 = 12) estimators agree with a total effect, the main effects, first-order interaction effects and second-order interaction effects.

	Single scores	Chi-squared values	Degrees of freedom
Total effect:	+27.59 * *	761.28	1
Main effects:			
A-effect:	+ 4.02 * *	16.16	1
B-effect:	-24.24 * *	587.35	1
C-effects:	- 5.98 * * }	265.27	2
	+14.19 * * }		
First-order interaction effects:			
A x B effect:	+ 0.41 N.S.	0.17	1
A x C effects:	+ 0.10 N.S. }	0.23	2
	- 0.46 N.S. }		
B x C effects:	- 4.04 * * }	43.26	2
	- 5.70 * * }		
Second-order interaction effects:			
A x B x C effects:	- 0.35 N.S. }	1.31	2
	+ 1.03 N.S. }		
N.B. 95% limit:	<u>±</u> 1.96	3.84	1
		5.99	2

Scores are in standard form. In this way it can already be ascertained which effects are significant, for instance at 5% level. These are marked * *. This test can also be made with an X^2 -test. For each effect we then find one X^2 -value (see column 2) with the relevant number of degrees of freedom (see column 3). Note that the X^2 -values for $df=1$ are equal to the square of the single scores.

In the Chi-squared tests it is then assumed in each case that all estimators agreeing with the effect are equal to zero, and the X^2 -value indicates the extent of the discrepancy between the model so obtained and the data.

Where we are concerned with one degree of freedom per effect, the significance of both tests is by definition identical. In this analysis, the other X^2 -values provide the same result as the single-score test. This is not necessarily always the case. The single scores, for instance, may all be (just) not significant, but together yield a significant X^2 -value.

It is also possible that only one single score is significant, which does not make the total X^2 -value significant.

In such cases the X^2 -value and the single scores thus provide additive information.

Interpretation of the data

In general the main effects and the total effect themselves are not very significant in interpreting the data. Here, however, where a correction was made for the number of inhabitants, it can be said as regards the A-effect that per inhabitant there are fewer accidents in the Rest of the Netherlands than in Noord-Brabant (the direction of the effect is shown by the sign!).

In order to interpret this phenomenon, we would have to know something for instance about the degree of urbanisation in Noord-Brabant and in the Rest of the Netherlands and also at least something about the number of traveller and vehicle kilometres. For interpretation of the B x C effect it is important to realise that this relates to cases where drinking was established. It would be

interesting to relate this to built-up/non-built-up areas as well. All this should make it clear that interpretation of the effects is an exercise unconnected with the analysis itself.