

-

-

1.00

Statistics Statistique Canada Canada

1.1



Methodology Branch

Census & Household Survey Methods Division

Direction de la méthodologie

Division des méthodes de recensement et d'enquêtes ménages





ANALYSIS OF CATEGORICAL DATA FROM SURVEYS WITH COMPLEX DESIGNS: SOME CANADIAN EXPERIENCES

D.A. Binder, M. Gratton, M.A. Hidiroglou, S. Kumar and J.N.K. Rao

Number: METH 85-083E

STATISTICS	-	ATISTIQUE				
CANADA	C	CANADA				
10 N	30	1992				
LIBRARY BIBLIOTHÈQUE						



Analysis of Categorical Data from Surveys with Complex Designs: Some Canadian Experiences¹

D.A. Binder, M. Gratton, M.A. Hidiroglou, S. Kumar and J.N.K. Rao²

ABSTRACT

Goodness of fit tests, tests for independence in a two-way contingency table, log-linear models and logistic regression models are investigated in the context of samples which are obtained from complex survey designs. Suggested approximations to the null distributions are reviewed and some examples from the Canada Health Survey and Canadian Labour Force Survey are given. Software implementation for using these methods is briefly discussed.

KEYWORDS: χ^2 statistic; Wald Statistics; Goodness of fit; Independence in two-way tables; Log-linear and logistic regression model.

1. INTRODUCTION

A sketch of the historical development of modern categorical data analysis has been given in the excellent review paper by Imrey, Koch and Stokes (1981). These techniques, applied in the context of random samples derived as independent selections from a common distribution function, are not directly applicable to survey samples collected using complex survey designs.

Koch *et al* (1975), Shuster and Downing (1976), developed asymptotically valid methods, based on the Wald statistic that take the survey design into account, but requiring access to the micro-data file or at least the full estimated covariance matrix of cell estimates. Cohen (1976) and Altham (1976) proposed a simple model for clustering and showed that the generalized Wald statistic for goodness of fit is a multiple of χ^2 , when the model holds. Brier (1978) considered a similar model, but studied general hypotheses on cell probabilities, and proved that a multiple of the corresponding Pearson statistic is asymptotically distributed as a χ^2 random variable, when the model holds. Fellegi (1980) deflated the χ^2 using a correction factor based on the mean of the estimated design effects. Fay (1985) developed jackknife χ^2 and G² statistics, also taking the design into account, but requiring the cell estimated at the primary sampling unit level. Rao and Scott (1981) developed a correction to χ^2 (or G²) based on the Satterthwaite to approximation to the asymptotic distribution of χ^2 , requiring the full estimated covariance matrix.

In this paper, we discuss the problems of fitting models and testing hypotheses with categorical data resulting from complex designs. For data collected using complex designs, some adjustments to the classical methods described by Imrey, Koch and Stokes (1981) are necessary in order to make valid inferences. If the published tables are provided along with the cell and marginal design

This paper is a revised and expanded version of that presented at the Seminar on Recent Developments in the Analysis of Large Scale Data Sets sponsored by Statistical Office of European Communities, November 16-18, 1983, Luxemburg,

D.A. Binder, Institutional & Agriculture Survey Methods Division, M. Gration, EDP Planning and Support Division, M.A. Hidiroglou, Business Survey Methods Division, S. Kumar, Census & Household Survey Methods Division, Statistics Canada, Tunney's Pasture, Ottawa, Ontario, Canada, K1A 016, and J.N.K. Rao, Department of Mathematics and Statistics, Carleton University, Ottawa, Ontario, CANADA.

144 WI - 12 Grammer 1

and the second sec

x - - 12

THE REPORT OF A DESCRIPTION OF A DESCRIP

All a billion and a state in the second seco

and the second state of th

and the second s

Binder et al.: Analysis of Categorical Data

effects, some of the approximations to the null distributions of our test statistics can be obtained, without having access to the complete micro-data file. On the other hand, for applications where the complete micro-data file is available, alternative approaches will be described.

For illustrative purposes, Section 2 begins with the standard goodness of fit problem. This discussion is then extended in Section 3 to tests of independence in a two-way contingency table. This leads to a general discussion of log-linear models in Section 4. Logistic regression models are described in Section 5. In Section 6 we summarize the existing situation with respect to software development for these methods at Statistics Canada. In Section 7, we discuss the appropriateness of these methods. Numerical examples are taken primarily from the Canada Health Survey. An application from the Canadian Labour Survey is given in Section 5.

2. GOODNESS OF FIT

2.1 Multinomial Sampling

Suppose we select n independent and identically distributed observations Y_1, \ldots, Y_n from a discrete distribution with k categories, where $\Pr(Y = i) = \pi_i$; $\sum_{i=1}^k \pi_i = 1$. We observe the random vector $\eta = (n_1, \ldots, n_{k-1})^T$, which has a multinomial distribution. Our estimate of $\pi = (\pi_1, \ldots, \pi_{k-1})^T$ is given by $p = \eta/n$. This estimate is unbiased and has covariance matrix given by $V\{p\} = (D_{\pi} - \pi\pi^T)/n = P/n$, where $D_{\pi} = \text{diag}\{\pi_1, \ldots, \pi_{k-1}\}$. Note that $P^{-1} = D_{\pi}^{-1} + (11^T/\pi_k)$. Asymptotically, $n^{i_1}(p - \pi) \to N(0, P)$. For a given π_o , the goodness of fit problem is to test the hypothesis.

 H_{c} : $\pi = \pi_{c}$,

against the alternative

$$H_1: \pi \neq \pi_0. \tag{2.1}$$

Letting P_o represent P evaluated at π_o , the Wald statistic for this test is

$$V_{i} = n(\underline{p} - \underline{\pi}_{o})^{T} P_{o}^{-1} (\underline{p} - \underline{\pi}_{o})$$
$$= n \sum_{i}^{K} \{ (p_{i} - \pi_{io})^{2} / \pi_{io} \},$$

which is the familiar Pearson chisquare test. Under H_o this is asymptotically χ^2_{k-1} . The likelihood ratio test for this problem is given by

$$LR_{i} = 2n \Sigma_{i} p_{i} \log(p_{i}/\pi_{i0}).$$

Since $2p_i \log(p_i/\pi_{io})$ is asymptotically equivalent to $2(p_i - \pi_{io}) + (p_i - \pi_{io})^2/\pi_{io}$ under H_o , we see that the likelihood ratio test is asymptotically equivalent to the Pearson chisquare statistic under H_o .

Another possible test for this hypothesis is derived by defining the vector of logs, $\underline{\mu}_{o} = \log \pi_{o}$ and $\hat{\underline{\mu}} = \log p$. Now under the null hypotheses $\hat{\underline{\mu}} - \underline{\mu}_{o}$ is asymptotically equivalent to $\underline{D}_{ro}^{-1}(\underline{p} - \underline{\pi}_{o})$. Therefore, $n^{\nu_{1}}(\hat{\underline{\mu}} - \underline{\mu}_{o}) \rightarrow N(\underline{0}, \underline{D}_{ro}^{-1} - \underline{1}\underline{1}^{T})$ under H_{o} and the Wald statistic is

$$\begin{split} W_{2} &= (\hat{\mu} - \mu_{o})^{T} \left[D_{\pi_{0}} + (\pi_{o} \pi_{o}^{T} / \pi_{k_{0}}) \right] (\hat{\mu} - \mu_{o}) \\ &= \sum_{i=1}^{k} \pi_{io} (\hat{\mu}_{i} - \mu_{io})^{2}, \end{split}$$

where $\mu_{ko} = \log \pi_{ko}$ and $\hat{\mu}_{k} = \log p_{k}$.



This approximation is obtained by noting that under H

$$\pi_{ko}(\hat{\mu}_{k} - \mu_{ko}) \doteq p_{k} - \pi_{ko}$$

= - $(p - \pi_{o})^{T} \underline{1} - (\mu - \mu_{o})^{T} \overline{z}_{1}$.

Note that W_{i} is also asymptotically equivalent to the Pearson chisquare test under H_{i} .

2.2 Other Sampling Schemes

These results for W_1 , W_2 and LR_1 are well-known. The question of interest to us here is the implication of the more general assumption that $n'(p-\pi) \rightarrow N(0, V)$, where I is not necessari-Iv equal to P. Here p is a survey estimate of π and may depend on sampling weights and other adjustment factors. This situation often arises in sampling under a complex sample design. We assume that \hat{V} is a consistent estimate of V. There are two approaches which we shall consider here. The first is to construct the appropriate Wald statistic for the given sample design. This would be

$$W_{3} = n(p - \pi_{0})^{T} \tilde{V}^{-1}(p - \pi_{0}),$$

where the rank of \hat{V} is k-1 so that W_3 is asymptotically χ^2_{k-1} under H_o . An alternative approach would be to use W_1 , W_2 or LR_1 directly as a test statistic. Now from multivariate normal theory, we know that the distribution of $n(p - \pi_p)^T P_p^{-1}(p - \pi_p)$ is that of $\sum \delta_i Z_i^2$, where $\{Z_i^2\}$ are independent χ_1^2 random variable and $\delta_i = (\delta_1, \ldots, \delta_{k-1})^T$ are the eigenvalues of $P_{o}^{-1}V$; see Johnson and Kotz (1970, pg. 150). This result was shown by Rao and Scott (1981), who call the δ_i 's generalized design effects. We note that for k = 2, we have $\delta = n\sigma_p^2/\{\pi_o(1 - \pi_o)\}\)$, where $\sigma_p^2 = V\{p\}$. This is the usual design effect for p under H_o .

2.3 Approximations

Now, in general, the distribution function for linear combinations of χ_1^2 random variables is cumbersome, although their moments are easily obtained. Rao and Scott (1981) have suggested two approximations to obtain the significance levels. The first is to approximate the distribution as being proportional to a χ^2_{k-1} random variable, the proportionality constant being determined by equating the mean of the approximating distribution to that of the theoretical distribution. This results in the approximation

$$\sum_{i=1}^{k-1} \delta_i Z_i^2 \doteq \left\{ \sum_{i=1}^{k-1} \delta_i / (k-1) \right\} \chi_{k-1}^2$$
(2.2)

Now.

$$\Sigma \delta_{i} = \operatorname{tr}(\underline{P}_{o}^{-1}\underline{V})$$
$$= \sum_{i=1}^{k} \underline{v}_{ii} / \pi_{io}$$
$$= \sum_{i=1}^{k} d_{i}(1 - \pi_{io}),$$

which depends only on the cell design effects $\{d_i\}$, where v_i is the *i*-th diagonal element of \hat{V} and $d_i = v_0 / [\pi_w (1 - \pi_w)]$. This approximation is particularly convenient when the full covariance matrix is not known, but the cell design effects are given. This is often the case for official published data.



Binder et al.: Analysis of Categorica Data

Census Age Distribution for Canada (1978-9)								
	Age							
Census	15-19	20-24	25-34	35-44	45-54	55-64	65+	Total
Distribution	.133	.127	.218	.152	.140	.115	.115	1.000
Distribution of those consuming 1-6 drinks/week	.117	.150	.264	.175	.148	.093	.053	1 000
Design Effect	1.4	1.2	2.2	1.1	0.6	1.1	1.0	

Table 1

Age Distribution Among Those Consuming 1-6 Drinks Per Week. Census Age Distribution for Canada (1978-9)

Example 1

For the Canada Health Survey (1978-9), a stratified multi-stage household survey, data was derived for the age distribution among those consuming one to six drinks per week, based on a sample of 5,204 persons, aged 15 years and over. A description of the survey may be found in "The Health of Canadians" (Statistics Canada Catalogue No. 82-538).

The data, taken from Hidiroglou and Rao (1981), are presented in Table 1. The raw value for W_1 is 298. This is reduced to 248 by taking the approximation given by (2.2). For these data, the post-stratification adjustments for age and sex lead to small design effects.

A second approximation to the distribution of $\sum \delta_i Z_i^2$, suggested by Rao and Scott (1981), is the Satterthwaite (1946) approximation: $\sum \delta_i Z_i^2 \approx a \chi_i^2$. To obtain a and ν , it is necessary to compute

Σ

$$\delta_i^2 = \operatorname{tr}\{(\underline{P}_o^{-1}\underline{V})^2\}$$

$$= \sum_{i=1}^{k} \sum_{j=1}^{k} \frac{v_{ij}^2}{(\pi_{i0}\pi_{j0})}.$$

However, this depends on all the terms of the matrix \hat{V} . The important point, though, is that some adjustment to the multinomial test statistic is necessary to obtain the appropriate significance level.

An alternative approximation, suggested by Fellegi (1980), is to divide the statistic $n(p - \pi_o)^T P_o^{-1}(p - \pi_o)$ by the average design effect. \vec{d} , instead of the weighted average given in (2.2). The effect of this on the data in Table 1 is that the adjusted chisquare value is 243, which is comparable to Rao and Scott's (1981) approximation.

3. TESTS OF INDEPENDENCE IN A TWO-WAY TABLE

3.1 Multinomial Sampling

We now suppose that the categories of the multinomial distribution can be cross-classified into an $r \times c$ table, where for the bivariate observation (Y_1, Y_2) we have $\Pr(Y_1 = i, Y_2 = j) = \pi_{ij}^*$ $\sum_{i=1}^r \sum_{j=1}^r \pi_{ij} = 1$. We denote $\pi_{i,2} = \sum_{j=1}^r \pi_{ij}$ and $\pi_{+j} = \sum_{i=1}^r \pi_{ij}^*$. We denote $\pi = (\pi_{1i}, \ldots, \pi_{1c}, \ldots, \pi_{c-1i}, \pi_{ij}, \ldots, \pi_{c-1i})^T$, $\underline{T}_R = D_{i_R} - \underline{\pi}_R \overline{\pi}_R^T$. $P_C = D_{i_R} - \underline{\pi}_C \overline{\pi}_C^T$. We observe the random vector \underline{n} from the multinomial distribution, where $E\{\underline{n}\} = n\underline{\pi}$. We let $\underline{p} = \underline{n}/n$, $p_{i_R} = \sum_{j=1}^r p_{i_R}^*$ and $p_{+j} = \sum_{j=1}^r p_{i_R}^*$.



We wish to test the hypothesis of independence

$$H_{i}: \pi_{i} - \pi_{i}, \pi_{i} = 0$$
 for $1 \le i \le r-1; 1 \le j \le c-1$.

against the alternative

$$H_i: \pi_i - \pi_i \pi_j \neq 0 \text{ for some } (i, j).$$

If we construct $h_{ij} = p_{ij} - p_{ij}p_{ij}$, for $1 \le i \le r-1$ and $1 \le j \le c-1$, then under multicomial sampling under H_{ij} , the asymptotic covariance matrix for $h_i = (h_{11}, \ldots, h_{1,c-1}, \ldots, h_{ij})^{-1}$ is $P_R \otimes P_c$, where \otimes denotes the direct matrix product operation. Hence, the Wald statistic under H_i becomes

$$W_{4} = h^{T} (\hat{P}_{C}^{-1} \otimes \hat{P}_{R}^{-1}) h$$
$$= \sum_{i=1}^{r} \sum_{j=1}^{c} (p_{ij} - p_{i,j} p_{j,j})^{2} / (p_{i,j} p_{j,j})$$

the familiar chisquare test with (r-1)(c-1) degrees of freedom.

Another test, which is asymptotically equivalent to W_4 under H_o , is the likelihood ratio test given by

$$LR_{2} = 2n \left[\sum_{i=1}^{c} \sum_{j=1}^{c} p_{ij} \log p_{ij} - \sum_{i=1}^{c} p_{i+1} \log p_{i+1} - \sum_{j=1}^{c} p_{+j} \log p_{-1} \right].$$

An alternative approach for this problem, which is a special case of the methods described in Grizzle, Starmer and Koch (1969) is to consider a Wald statistic based on

$$\{f_n = \log p_n - \log p_1 - \log p_1; \text{ for } 1 \le i \le r-1, \text{ and } 1 \le j \le c-1\}$$

The asymptotic covariance matrix for $\hat{f} = (f_{11}, \ldots, f_{l,c-1}, \ldots, f_{r-1,c-1})^T$ is $(D_{r_c}^{-1} - 1 1^T)$ $\otimes (D_{r_c}^{-1} - 11^T)$. Therefore the Wald statistic becomes

$$W_{5} = \int_{-}^{T} \left[\left(\hat{D}_{\pi_{R}} + \frac{\pi_{R} \pi_{R}}{p_{r+1}} \right) \otimes \left(\hat{D}_{\pi_{C}} + \frac{\pi_{C} \pi_{C}}{p_{r+1}} \right) \right] f$$

Now under H_o we note that f_{ij} is asymptotically equivalent to

$$\frac{p_{ij}}{\pi_{i+}\pi_{+j}} - \frac{p_{i+}}{\pi_{i+}} - \frac{p_{+j}}{\pi_{+j}} + 1,$$

so that $\sum_{i=1}^{r} \pi_{i+} f_{ij} = \sum_{j=1}^{c} \pi_{+j} f_{ij} = 0$. Using this approximation W_{s} becomes

$$W'_{5} = \sum_{i=1}^{r} \sum_{j=1}^{c} p_{i+} p_{+j} f_{ij}^{2}$$

It should be noted that under H_o , the statistics W_a , LR_c and W_c are all asymptotically equivalent to

$$\sum_{i=1}^{r} \sum_{j=1}^{c} \frac{(p_{ij} - \pi_{i+}\pi_{+j})^2}{\pi_{i+}\pi_{+j}} = \sum_{i=1}^{r} \frac{(p_{i+} - \pi_{i+})^2}{\pi_{i+}} = \sum_{j=1}^{c} \frac{(p_{+j} - \pi_{+j})^2}{\pi_{+j}}$$
(3.1)

This result will prove useful in Section 3.3.



3.2 Other Sampling Schemes

Now, relaxing the assumption that n is multinomial, we assume instead that $n'(p-\pi) \rightarrow N(Q, V)$ where p is a survey estimate which may depend on sampling weights and other adjustment factors. For this case, Shuster and Downing (1976) and Fellegi (1980) suggest that we construct the Wald statistic based on $\{h_{ij} = p_{ij} - p_{ij}, p_{jj}\}$. If we let J_a be the $(a-1) \times a$ matrix given by

$$J_{a} = \begin{bmatrix} I & Q \end{bmatrix}$$
(3.2)

and let $\hat{H} = (J_r - \pi_R 1^T) \otimes (J_c - \pi_C 1^T) - (\pi_R 1^T \otimes \pi_C 1^T)$

then the Wald statistics is

$$W_{\lambda} = h^{\gamma} (\hat{H} \hat{V} \hat{H}^{\gamma})^{-1} h,$$

which under H_0 is asymptotically $\chi^2_{(r-1)(r-1)}$.

Alternatively, we could construct a Wald statistic based on $\{f_{ij} = \log p_{ij} - \log p_{ij} - \log p_{ij}\}$. This is a special case of the log-linear model approach to be discussed in Section 4. We define $(r - 1) \times r$ and $(c - 1) \times c$ matrices as follows:

$$\hat{E}_{\mathsf{R}} = \left[\hat{D}_{\pi_s}^{-1} \mid \mathbf{Q} \right], \quad \hat{E}_{C} = \left[\hat{D}_{\pi_c}^{-1} \mid \mathbf{Q} \right].$$

We let $\hat{F} = (\hat{E}_R - \hat{E}_R \mathbb{1} \mathbb{1}^T) \otimes (\hat{E}_C - \hat{E}_C \mathbb{1} \mathbb{1}^T) - (\hat{E}_R \mathbb{1} \mathbb{1}^T \otimes \hat{E}_C \mathbb{1} \mathbb{1}^T)$.

The appropriate Wald statistics is

$$W_{\tau} = f^{T} (\hat{F} \hat{V} \hat{F}^{T})^{-1} f.$$

Now, analogously to the goodness of fit problem in Section 2, Rao and Scott (1981) have considered null distributions of the test statistics based on W_4 , LR_2 and W_5 , which are all asymptotically equivalent to the null distribution of (3.1). We see, therefore, that the null distribution is the same as $\sum_{i=1}^{(r-1)(c-1)} \delta_i Z_i^2$, where $\{Z_i^2\}$ are independent χ_1^2 and the δ_i 's are the eigenvalues of

$$(P_{P}^{-1} \otimes P_{C}^{-1})(HVH^{T})$$
.

Cowan and Binder (1978) investigated the properties of the eigenvalue from a simple two-stage self-weighting design for a 2×2 table. They found that the eigenvalue increases as the degree of independence of the cell proportions within the primary sampling units decreased.



3.3 Approximations

An approximation for the distribution of $\Sigma \delta_i Z_i^2$ is

$$\sum \delta_i Z_i^2 \approx \frac{\sum \delta_i}{(r-1)(c-1)} \chi^2_{(r-1)(c-1)},$$

as in (2.2). Since the statistic is asymptotically equivalent to(3.1) under H_o , by computing the mean of (3.1) we obtain

$$\sum \delta_{i} = \sum_{i=1}^{r} \sum_{j=1}^{c} d_{ij} \left(1 - \pi_{i+} \pi_{+j} \right) - \sum_{i=1}^{r} d_{i}^{(r)} \left(1 - \pi_{i-} \right) - \sum_{j=1}^{c} d_{j}^{(c)} \left(1 - \pi_{-j} \right),$$

where d_{ij} is the cell design effect; $d_{ij}^{(i)}$ and $d_{ij}^{(i)}$ are the row and column margin design effects, respectively. This particularly simple expression was obtained by Rao and Scott (1983). Fellegi (1980) suggested an alternative approximation as:

$$\left(\sum_{i=1}^{r}\sum_{j=1}^{c}d_{ij}/rc\right)\chi^{2}(r-1)(c-1)$$

Example 2

In Table 2, we give a 4 × 2 table from the Canada Health Survey, which cross-classifies drug use (four categories; 0, 1, 2, 3 + drug classes in a 2-day period) and sex (male, female). Here n = 31,668.

The raw value for W_4 is 774. Rao and Scott's (1981) adjustment reduces this to 437. Fellegi's (1980) adjustment reduces this to 327. The Wald statistics, W_6 , is 538. Hidiroglou and Rao (1981) found that the Rao and Scott (1981) approximation performs quite well relative to the Satterthwaite (1946) approximation which is based on the complete covariance matrix.

LOG-LINEAR MODELS

4.1 Multinomial Sampling

We now extend the results of the previous section to more general cross-classifications of the multinomial distribution. The standard results for these models are given in Bishop, Fienberg

Variety of Drugs Taken by Sex for Canada (1978-79)						
Sex			Numb			
		0	1	2	3 ±	Total
Male	Proportion Design Effect	0.293 1.56	0.134 3.37	0.048	0.021 1.38	0.496 0.00*
Female	Proportion Design Effect	0.228 3.59	0.159 3.13	0.072 2.85	0.045 1.96	0.504
Total	Proportion Design Effect	0.521 6.03	0.293 6.46	0.120	0.066 2.57	1.000

Table 2

* Because of age-sex post-stratification, these design effects are zero.



Binder et al.: Analysis of Categorical Data

and Holland (1975) and Fienberg (1980). We have $\pi = (\pi_1, ..., \pi_k)^r$ is a vector of cell proportions; $\sum_{i=1}^k \pi_k = 1$. We observe $\eta = (n_1, ..., n_k)^r$, the counts in each cell from a random sample, so that η has a multinomial distribution ($\sum n_i = n$). We let $p = \pi/n$ and define

 $\mu = \log \pi$.

The log-linear model assumes that for a parameter vector $\theta = (\theta_1, \ldots, \theta_n)^T$, we have

$$\mu(\theta) = u(\theta) \mathbf{1} + X\theta,$$

where X is a known $k \times t$ matrix of full rank and $X^{T} = 0$. Note that $t \leq k-1$. If t = k-1, we have the saturated model.

The maximum likelihood estimate for θ is given by solving

$$X^{T}(p - \hat{\pi}) = 0,$$
 (4.1)

where $\hat{\pi} = \pi(\hat{\theta})$. Now, asymptotically we have

$$\hat{\pi} - \pi \doteq PX(\theta - \theta) ,$$

where $P = D_{\tau} - \pi \pi^{T}$. From (4.1), we then obtain

$$\hat{\theta} - \theta \doteq (X^T P X)^{-1} X^T (p - \pi)$$

and

$$\tau - \pi \doteq P X (X^T P X)^{-1} X^T (p - \pi).$$

Since $n^{1/2}(p-\pi) \rightarrow N(0, P)$ we obtain

$$n^{\frac{1}{2}} (\hat{\theta} - \theta) \rightarrow N [0, (X^T P X)^{-1}]$$
$$n^{\frac{1}{2}} (\pi - \pi) \rightarrow N [0, P X (X^T P X)^{-1} X^T P]$$

Suppose now that the linear expression $X\theta$ can be decomposed as $X_1\theta_1 + X_2\theta_2$ where X_1 and X_2 are full rank, X_1 is $k \times r$, X_2 is $k \times s$, θ_1 is $r \times I$ and θ_2 is $s \times I$, where r + s = t. We consider the problem of testing

$$H_{0}:\theta_{1}=0$$

against the alternative

$$H_1: \theta_2 \neq 0.$$

We use θ_1 , θ_2 , π , etc. to denote the estimates under the full model H_1 . Alternatively, we let $\hat{\theta}_1$, $\hat{\pi}$, to denote estimates under H_o .

Now,

$$n^{1/2}(\theta_{1}, -\theta_{2}) \rightarrow N[0, (\tilde{X}_{1}^{T}P\tilde{X}_{2})^{-1}]$$

where

$$\tilde{X}_{2} = \begin{bmatrix} I - X_{1} (X_{1}^{T} P X_{1})^{-1} X_{1}^{T} P \end{bmatrix} X_{2}$$
(4.2)



so that the Wald statistic is

$$W_{\rm g} = n\hat{\theta}_{2}^{T}\hat{X}_{2}^{T}\hat{P}\hat{X}_{2}\theta$$

Under H_{o} , this is asymptotically equivalent to the Pearson chisquare statistic

$$n\left(\hat{\pi} - \hat{\pi}\right)^T \hat{D}^{-1}\left(\hat{\pi} - \hat{\pi}\right),$$

or the likelihood ratio test

$$LR_{1} = 2n \hat{\Sigma}_{1} p \log(\hat{\pi} / \hat{\pi}).$$

Under H_{p} , these statistics are asymptotically χ_{1}^{2} .

4.2 Other Sampling Schemes

We still assume that the cell proportions, π , satisfy $\mu = \log \pi - u(\theta_1, \theta_2) [+ X_1 \theta_1 + X_2, \theta_2]$, but we now have $n^{(n)}(p - \pi) \rightarrow N(0, V)$, where p is a survey estimate.

Rao and Scott (1983) suggest the following Wald statistic for testing $\theta_2 = 0$. We let ζ be any $k \times s$ matrix with $\zeta^T X_1 = 0$, $\zeta^T 1 = 0$ and $\zeta^T X_2$ nonsingular. For example if $X_1^T X_2 = 0$ then $\zeta = X_2$ is convenient. Now the hypothesis is equivalent to $\zeta^T \mu = 0$. We have

$$C^{T}(\hat{\mu} - \mu) \doteq C^{T} D^{-1}(\hat{\pi} - \pi)$$

$$= C^T X (X^T P X)^{-1} X^T (p - \pi),$$

where π is obtained from (4.1), based on the survey estimate, p.

We therefore have the Wald statistics

$$W_{o} = n\hat{\mu}^{T}C[C^{T}X(X^{T}\hat{P}X)^{-1}(X^{T}\hat{V}X)(X^{T}\hat{P}X)^{-1}X^{T}C]^{-1}C^{T}\hat{\mu}.$$

Similar results were also given in Binder (1983). If under H_1 , the model is saturated (r+s = k-1), then $p = \pi$ and we obtain

$$W_{o} = n\hat{\mu}^{T} C \left[C^{T} \hat{D}^{-1} \hat{V} \hat{D}^{-1} C \right]^{-1} C^{T} \hat{\mu}.$$

Rao and Scott (1984) show that if we use \hat{P} instead of \hat{V} in W_g then these are asymptotically equivalent to the likelihood ratio or Pearson χ^2 test statistics. They also show that the likelihood ratio test statistics is distributed as $\sum_{i=1}^{3} \delta_i Z_i^2$ under H_o , where $\{Z_i^2\}$ are independent χ^2 and $\{\delta_i\}$ are the eigenvalues of

$$(\tilde{X}_{\cdot}^{T}P\tilde{X}_{\cdot})^{-1}(\tilde{X}_{\cdot}^{T}V\tilde{X}_{\cdot}), \qquad (4.3)$$

for \hat{X} , defined in (4.2).

4.3 Approximations

As before, we approximate the null distribution

$$\sum_{i=1}^{s} \delta_{i} Z_{i}^{2} \approx \left(\frac{\sum \delta_{i}}{S}\right) \chi_{s}^{2}.$$

This involves computing the trace of (4.3). Rao and Scott (1984) show that if the model admits explicit solutions for both $\hat{\pi}$ and $\hat{\pi}$, then the approximation depends on the matrix V only through cell design effects and marginal design effects. This observation is particularly convenient when ony the estimated design effects for the cell proportions and margins are available, as is often the case for published tables.



Example 3

Hidiroglou and Rao (1983) considered all direct estimates from the three-way table: Drug use (5 categories: 0, 1, 2, 3, 4 + drug classes in a 2 day period) × Age (4 categories; 0-14, 15-44, 45-64, 65 +) × Sex (male, female), taken from the Canada health Survey. We give the results for testing whether Age and Sex are independent in each drug category (n = 31,668). This is equivalent to the hypothesis

$$H_{a}: \pi_{ak} = \pi_{ak} - \frac{\pi_{i+k}}{\pi_{i+k}} \cdot$$

Using Bishop, Fienberg and Holland's (1975) notation, where $\log \pi_{ijk} = u + u_{1ij} + u_{2ij}$ + $u_{1ijk} + u_{12ijk} + u_{12ijk} + u_{23ijk} + u_{123ijk}$, the null hypothesis is equivalent to

$$H_0: u_{23(k)} = u_{123(k)} = 0$$
 for all (i, j, k) .

The raw chisquare value is 23 based on 15 degrees of freedom. The average eigenvalue is 1.39, so that the approximation reduces the chisquare value to 16. Whereas the unadjusted chisquare value would lead the analyst to reject the hypothesis at the 10% level, the approximation indicates that $h_{\rm o}$ cannot be rejected even at the 30% level.

5. LOGISTIC REGRESSION MODELS

5.1 Multinomial Sampling

We now consider a logistic regression model for the conditional distribution of a binary response variable y given the vector \underline{x} of independent variables. In particular, this conditional distribution is

$$\Pr(y_i \mid x_i) = \pi(x_i)^{y_i} \left[1 - \pi(x_i)\right]^{1-y_i},$$

where $v \in \{0, 1\}$.

For the logistic regression model, we have

$$\log \left\{ \frac{\pi(\underline{x}_i)}{1 - \pi(\underline{x}_i)} \right\} = \underline{x}_i^T \underline{\theta},$$

where e is an unknown vector of parameters.

We note that if x_i is a categorical vector of 0's and 1's, this is a special case of a log-linear model as described in Section 4. Here we allow x_i to be arbitrary. The extension to the case of k-categories for the y-variable is straight-forward, it is also possible to generalize the model to

$$\log \left\{ \frac{\pi(\underline{x}_i)}{1 - \pi(\underline{x}_i)} \right\} = f(\underline{x}_i^T \theta),$$

for a known function f(), but we do not discuss this here.

Now, the maximum likelihood estimate for θ is given by

$$X'(y - \hat{\pi}) = 0$$

where $y = (y_1, \ldots, y_n)^T$, $\tilde{\pi} = [\hat{\pi}(y_1), \ldots, \hat{\pi}(x_n)]^T$ and $X = [x_1, \ldots, x_n]^T$. Under suitable regularity conditions, we have

$$n(\theta - \theta) \rightarrow N[0, n(X^{T}\Lambda X)^{-1}], \text{ where } \Lambda = D(1 - D).$$



If we have $X\theta = X_1\theta_1 + X_2\theta_2$ and consider testing the hypothesis

$$H_0: \theta_2 = 0$$
$$H_1: \theta_2 \neq 0,$$

we obtain the Wald statistic

$$W_{10} = n \hat{\theta}_{1}^{T} (\hat{X}_{1} \wedge X_{1}) \hat{\theta}_{2}$$

where

$$X_{1} = \begin{bmatrix} I & -X_{1}(X_{1}^{T}\Lambda X_{1})^{-1}X_{1}^{T}\Lambda \end{bmatrix} X_{2}$$

The likelihood ratio test here is

$$LR_{4} = 2\sum_{i=1}^{n} \left[y_{i} \log \left(\frac{\hat{\pi}_{i}}{\hat{\pi}_{i}} \right) + (1 - y_{i}) \log \left\{ \frac{(1 - \hat{\pi}_{i})}{(1 - \hat{\pi}_{i})} \right\} \right]$$

which is asymptotically equivalent to W_{10} under H_0 .

5.2 Other Sampling Schemes

Suppose now that $n^{-\frac{1}{2}}X^T(y - \pi) \to N(0, V)$ and that \hat{V} is a consistent estimator of V. Here y is not necessarily a vector of 0's and 1's, but may in fact depend on the sampling weights and other adjustment factors. Estimating V is usually possible since $X^T(y - \pi)$ is the sum of random observations and most sample designs admit a consistent estimator of the sum of (not necessarily independent) observations. To estimate V we use $\hat{\pi}$ instead of π in the estimate. Since asymptotically

$$(\hat{\theta} - \theta) \doteq (X^T \Lambda X)^{-1} X^T (y - \pi),$$

we have that

$$n^{\nu_1}(\hat{\theta} - \theta) \rightarrow N[0, n^2(X^T \Lambda X)^{-1} V(X^T \Lambda X)^T];$$

see Binder (1983) for a detailed justification of this result. Now, a Wald statistic may be constructed from the estimated covariance matrix for θ_2 .

Table 3

Logistic Regression Model for Explaining Use of Physician Services

Variable	Туре	d.f.	Wald Statistic	
Age	Categorical	4	19.232	
Sex	Categorical	1	12.494	
Age-Sex Interactions	Categorical	4	36.001	
Family Income	Categorical	5	14.642	
Occupation	Calegorical	3	8.614	
Occupation-Sex Interactions	Categorical	3	11.501	
Marital Status	Categorical	3	45.752	
Medical History	Categorical	2	36.700	
Number of Health Problems	Quantitative	1	81.554	
Drug Use	Categorical	2	272.175	
Number of Accidents	Quantitative	2	106.372	
Number of Disability Days	Quantitative	2	29.052	
Community Size	Categorical	2	11.751	
Provincial Physician -				
Population Ratios	Quantitative	1	0.540	

Example 4

A logistic regression model was fit on 20,726 respondents from the Canada health survey to explain use or non-use of physician services over a 12-month period. In total it was estimated that 77% of the population visited a physician at least once. The results are summarized in Table 3. For more complete details, see Binder (1983). The logistic model seemed to fit the data very well.

5.3 Qualitative Explanatory Variables

The theory of this section was obtained by G. Roberts in an unpublished manuscript (Carleton University). Here the explanatory variables are all qualitative. We label the domains, $\{1, \ldots, I\}$. We let p_i be the survey estimate of the *i*-th domain proportion and \hat{N}_i is the estimate of the size of the *i*-th domain, N_i . Under the model, the expected proportion in the *i*-th domain is f_i , where

$$\log \left\{ f/(l-f) \right\} = a^{T}\theta,$$

for \underline{a}_i known and $\underline{\theta}_i$ an unknown parameter. We define $\underline{A} = [\underline{a}_1, \ldots, \underline{a}_l]^T$ and let $D_{ij} = \text{diag} \{\hat{N}_{ij}, \ldots, \hat{N}_l\}$.

Under the model, the survey estimator of $f = (f_1, \ldots, f_l)^T$ is given by \hat{f} , the solution to

$$A^{T}D_{S}(p-\hat{f}_{s}) = 0.$$
 (5.1)

Since asymptotically

$$\hat{\theta} - \theta \doteq (A^T \Delta A)^{-1} A^T D_{\lambda} (p-f),$$

where $\Delta = \text{diag}\{N_1f_1(1-f_1), \ldots, N_lf_l(1-f_l)\}$, we have

$$n^{\nu_1}(\theta - \theta) \rightarrow N[0, (A^T \Delta A)^{-1} A^T D, V, D, A (A^T \Delta A)^{-1}]$$

whenever $n^{1/2}(p-f) \rightarrow N(0, V_p)$.

Under independent binomial sampling, the covariance matrix reduces to $(N/n)(A^T \Delta A)^{-1}$, where n is the sample size.

The likelihood ratio test for testing goodness of fit is

$$LR_{5} = 2(n/\hat{N}) \sum_{i=1}^{l} \hat{N}_{i} [p_{i} \log(p_{i}/\hat{f}_{i}) + (1-p_{i}) \log\{(1-p_{i})/(1-\hat{f}_{i})\}],$$

where n is the sample size and $\hat{N} = \sum \hat{N}_{r}$. Under H_{r} this is asymptotically equivalent to

$$W_{11} = (n/\hat{N}) \sum_{i=1}^{l} \hat{N}_{i} (p_{i} - \hat{f}_{i})^{2} / [f_{i} (1 - f_{i})].$$

In general, the distribution of LR_s will be that of $\Sigma \delta_i Z_i^2$, where $\{Z_i\}$ are independent χ_1^2 , and $\{\delta_i\}$ are the eigenvalues of $\hat{N}^{-1} D_{\hat{N}} [\Delta^{-1} - A(A^T \Delta A)^{-1} A^T] D_{\hat{N}} V_p D_{\hat{N}} [\Delta^{-1} - A(A^T \Delta A)^{-1} A^T] \Delta D_{\hat{N}}^{-1}$. By taking the expectation of W_{11} , and approximating

$$W_{11} \approx \frac{\sum \delta_i}{I-s} \chi^2_{I-s}$$

where $s = \operatorname{rank}(A)$, we obtain

$$\Sigma \delta_i = (n/\hat{N}) \sum_{i=1}^{l} \hat{N}_i v_n^{(n)} \{ f_i (1-f_i) \}$$

where $v_{ii}^{(i)} = V \{ p_i - \hat{f}_i \}$. The $\{ v_{ii}^{(i)} \}$ may be computed using the relationship $p_i - \hat{f}_i \doteq [I - \text{diag} \{ f_i(1 - f_i) \} A (A^T \Delta A)^{-1} A^T D_N] (p_i - f_i).$

Example 5

The data from the October 1980 Canadian Labour Force Survey was used to fit logistic (logit) models for the probability of being employed. The sample consisted of males aged 15-64 who were in the labour force and not full time students. A logit model, quadratic in age and in education, was fitted. Age-group levels were formed by dividing the interval [15, 64] into ten groups with the jth age-group being the interval [10 + 5j, 14 + 5j], j = 1, 2, ..., 10. The midpoint of each age-group was used as the value of the age for all persons in that age-group. Six levels of education were formed by assigning to each person a value based on the median years of schooling. Age by education classification led to the formation of 60 cells.

Let $\pi_i = \Pr\{\text{an individual in the ith cell is employed}\}, i = 1, 2, ..., 60$. We assume that $0 < \pi_i < 1$. Hence $1 - \pi_i$ represents the probability that the individual in the ith cell is unemployed. The model, considered for fit, was

$$\ln \frac{\pi_i}{1 - \pi_i} = \beta_0 + \beta_1 a_i + \beta_2 a_i^2 + \beta_3 d_i + \beta_4 d_i^2, \qquad (1)$$
$$i = 1, 2, \dots, 60$$

where a_i and d_j are the age and education variable values for the individuals in the ith cell.

Using the survey estimates p_i of π_p the values of Pearson's statistic W_{11} and the likelihood ratio statistic LR_s were computed as $W_{11} = 98.94$ and $LR_s = 101.20$. The upper 5% point of the chi-square distribution, with 55 degrees of freedom, is 73.31. Using these values of W_{11} or LR_s we would reject the model I. These values of W_{11} or LR_s , however, are appropriate only if the sample was a random sample.

The estimate average eigenvalue, $\sum \delta_i/55$, for testing goodness of fit for this data is 1.88. This would reduce W_{11} to 52.63 and LR_5 to 53.83. Hence, with this adjustment, we find that the data are consistent with the model (1).

The use of the Wald statistic, $(p - \hat{f})^{T} [\hat{V}^{(r)}]^{-} (p - \hat{f})$, for testing the goodness of fit was also considered. Here we use the g-inverse of $\hat{V}^{(r)}$ since the matrix is singular. Some perturbation to the estimates of p_i , when $p_i = 1$, was necessary for computing the Wald statistic. It was found that the Wald statistic was unstable for our problem. Minor perturbations in the estimates of p led to considerable change in the value of the Wald statistic. Also the value of the Wald statistic is very large here due to instability in the estimated covariance matrix involved in its calculation. The Wald statistic is at least 30 times larger than our adjusted Chi-squared values.

6. SOFTWARE CONSIDERATIONS

Advancement of computer technology has made data collection, storage and retrieval operations easy and efficient. Powerful generalized software systems, such as TPL, STATPAK and ESTIMATION SYSTEM, have been used to produce cell estimates and some of their variances fairly easily to users and analysts. As well a number of commercially available packages such as BMDP, SPSS and SAS are powerful analytic tools in certain contexts. However, the ability to perform analysis such as those described in this paper are limited. For example, in situations involving hypothesis testing or statistical inference, these packages assume that the data to be analyzed come from surveys with simple random samples.

At present, an integrated software package, similar to the ones mentioned above, but designed for analyses of the type of data discussed in this paper, is not available. As a result, the researcher requiring a quick solution to his problem is usually forced to use existing statistical packages which may not be appropriate.

Binder et al. Analysis of Categorical Data

The alternatives are

- use existing packages with modifications
- use existing stand-alone software
- write customized programs
- use combinations of the above.

For the analyses given in this paper, modifications to the MINI CARP program (Hidiroglou, Fuller and Hickman; 1980) were incorporated to obtain the results in Examples, 1, 2 and 3. For Example 4, a combination of PL/1 and SAS programs were developed. The analysis of the Labour Force Survey data (Example 5) used a combination of customized programs and SAS.

For the above alternatives, some practical drawbacks have been experienced, they include:

- (a) If an existing package is to be modified, intimate knowledge of the package is often required;
- (b) Identical information may have to be duplicated on separate data files, as these alternatives are not integrable like generalized systems;
- (c) Compared to an integrated "user-friendly" package, these alternatives lack elegance and operational efficiency as software;
- (d) Comprehensive documentation is not generally available for specially written programs limiting the availability of software.

Work is now ongoing to develop SAS based procedures for performing many of these analyses. Our ultimate goal is similar to that proposed by Shah (1981); namely, the development of an integrated software package for survey data analysis. This is a goal worth striving for, if we are to avoid the frustrations now being experienced by researchers who are faced with either developing their own software or using existing software which could lead to erroneous results and conslusions.

7. DISCUSSION

We have examined a number of problems which arise when fitting models to categorical data which have been collected under complex sampling designs. The basic approach has been to derive the appropriate Wald statistic for the fitted model or to use the test statistic which is motivated from multinomial-type sampling designs and find a suitable approximation to its null distribution.

We have not addressed the issue as to whether one should really be taking a model-based or design-based approach to begin with. Instead, we have concentrated on design-based inferences.

To put this issue into focus, let us reconsider the test of independence in a two-way contingency table. The question of independence arises if we are interested in whether knowing the value of variable Y_1 affects our knowledge about variable Y_2 . If it does not, for all the individuals in the population, then we say the variables are independent. However, if we also know the value of Y_3 , it may turn out that Y_1 and Y_2 are no longer independent. This is particularly important when Y_3 is a design variable (such as geographic stratum). Since design variables are usually known for all sampled individuals, we have one of two options: (a) we can say that the question of independence is no longer relevant, or (b) we can marginalize out Y_3 , and say that we are only interested in Y_1 and Y_2 , unconditionally. Assuming that we take approach (b), the results of this paper seem appropriate. In some cases it may be possible to test if Y_1 and Y_2 are conditionally independent given Y_3 .

There is a further difficulty, however. Suppose we are interested in the cell proportions π_{ij} from a finite population of size N. If we were to take a census from this population, it is highly unlikely that we would obtain $\pi_{ij} = \pi_{i+} \pi_{+j}$ exactly. The best that we could hope for is that some measure of association such as $N \Sigma \Sigma (\pi_{ij} - \pi_{i+} \pi_{+j})^2 / \pi_{i+} \pi_{+j}$ is small. Note that even

under a super-population model of exact independence, we would not expect this measure of association to be zero. Perhaps, we should instead be testing hypotheses such as

*H*_.: Measure of Association $\leq C$

H: Measure of Association > C.

Further research is needed in this area. However, for practical circumstances where the sampling fraction is not large, the methods given in this paper are suitable.

REFERENCES

- ALTHAM, P.A.E. (1976). Discrete vriable analysis for idividuals grouped into families. *Biometrika*, 63, 263-269.
- BINDER D.A. (1983). On the variances of asymptotically normal estimators from complex surveys. International Statistical Review, 51, 279-292.
- BISHOP, Y.M.M., FIENBERG, S.E. and HOLLAND, P.W. (1975). Discrete Multivariate Analysis: Theory and Practice. Cambridge, Massachusetts: MIT Press.
- BRIER, S.S. (1978). Discrete Data Models with Random Effects. Technical Report, University of Minnesota, School of Statistics.
- COHEN, J.E. (1976). The distribution of the chi-squared statistic under cluster sampling from contingency tables. Journal of the American Statistical Association, 71, 665-670.
- COWAN, J. and BINDER, D.A. (1978). The effect of a two-stage sample design on tests of independence in a 2 by 2 table. Survey Methodology, 4, 16-28.
- FAY, R.E. (1985). A jackknifed chi-squared test for complex samples. Journal of American Statistical Association, 80, 148-157.
- FELLEGI, I.P. (1980). Approximate tests of independence and goodness of fit based on stratified multistage samples. Journal of American Statistical Assocociation, 71, 665-670.
- FIENBERG, S.E. (1980). The Analysis of Cross Classified Data, (2nd ed.). Cambridge, Massachusetts: MIT Press.
- GRIZZEL, J.E., STARMER, C.F. and KOCH, G.G. (1969). Analysis of categorical data by linear models. Biometrics, 25, 489-504.
- HIDIROGLOU, M.A., FULLER, W.A. and HICKMAN, R.D. (1980). MINICARP: A program for estimating simple descriptive statistics and their variances for multi-stage stratified designs. Iowa State University: Ames, Iowa.
- HIDIROGLOU, M.A. and RAO, J.N.K. (1981). Chisquare tests for the analysis of categorical data from the Canada Health Survey. Paper presented at the International Statistical Institute Meetings, Buenos Aires, 1981.
- HIDIROGLOU, M.A. and RAO, J.N.K. (1983). Chi-square tests for the analysis of three-way contingency tables from the Canada Health Survey. Technical Report, Statistics Canada.
- IMREY, P.B., KOCH, G.G. and STOKES (1981). Categorical data analysis: Some reflections on the log linear model and logistic regression. part I: Historical and Methodological Ovewview. International Statistical Review, 49, 265-283.
- JOHNSON, N.L. and KOTZ, S. (1970). Continuous Univariate Distributions. Boston: Houghton Molflin.
- KOCH, G.G., FREEMAN, D.H. JR., and FREEMAN, J.L. (1975). Strategies in the Multivariate Analysis of Data from Complex Surveys. International Statistical Review, 43, 59-78.
- RAO, J.N.K. and SCOTT, A.J. (1981). The analysis of categorical data from complex sample surveys: Chi-square tests for goodness of fit and independence in two-way tables. *Journal of American Statistical* Association, 76, 221-30.
- RAO, J.N.K. and SCOTT, A.J. (1984). On chi-squared tests for multiway contingency tables with cell proportions estimated from survey data. *The Annals of Statistics*, 12, 48-60.

Binder et al.: Analysis of Categorical Data

- SATTERTHWAITE, F.E. (1946). An approximate distribution of estimates of variance components. Biometrics, 2, 110-114.
- SHAH, B.V. (1981). Development of survey data analysis software. Research Triangle Institute, Research Triangle Park, North Carolina.
- SHUSTER, J.J. and DOWNING, D.J. (1976). Two-way contingency tables for complex sampling schemes. Biometrika, 63, 271-276.

Ca OOS

