

THE ANALYSIS OF FREQUENCY DATA BASED ON AN ALTERNATIVE FORM  
OF CONSTRUCTING A CONTINGENCY TABLE

Jeffrey R. Wilson, Arizona State University  
William Warde, Oklahoma State University

SUMMARY

The effects of obtaining a test statistic for the test of independence when the data are obtained from a certain common complex sampling scheme is examined in this paper. The data are summarized in a two way contingency table where the cell entries are not necessarily an integer value. This alternative method of obtaining a two way table based on ratio estimators is compared to the traditional summing procedure in the construction of contingency tables. A Wald test statistic based on Wald (1943) for the test of independence is obtained for this particular sampling scheme and compared to these two forms of table construction techniques. This paper shows that an alternative approximation to the Wald test statistic for independence is to construct a Pearson type statistic based on the alternative table presented here rather than constructing a Pearson Statistic on the usual contingency table, as is sometimes done in complex sampling schemes.

1. Introduction

Test statistics are obtained for the tests of homogeneity and independence under a stratified cluster sampling scheme. In the past, several researchers have resorted to the use of the Pearson statistic in the presence of complex sampling procedures. Rao and Scott (1981, 1984), examined the behavior of the  $X^2$  statistic under complex designs by examining the eigenvalues of the product of the inverse of the covariance matrix under simple random sampling and the covariance matrix for the actual sampling scheme. Holt, Scott, and Ewings (1980) showed that a correction factor for  $X^2$  based on the design effects works well for the test of homogeneity. However, they demonstrated that for tests of independence an appropriate modifying factor is more difficult to complete. Cohen (1976), Brier (1980), Wilson and Koehler (1984), Wilson (1984) have considered models as a means of reducing the sample size required for variance estimation and producing useful test statistics.

In this paper consideration is given to an alternative form of the construction of the contingency table under the specified sampling scheme. The researcher is advised in the summarizing of the data. Test statistics are constructed based on the traditional form of table construction and based on the alternative form of construction. The covariance matrix for each of these cases is constructed based on the assumption that the cell proportions are multinomially distributed. A Wald test statistic is obtained based on the actual design and ignoring the multinomial assumption. A comparison is made with those

statistics obtained under the multinomial assumption and the Wald test statistic. A numerical example based on data obtained from a Wild life study Rolley and Warde (1985) is given in section 6 to demonstrate some of the results.

2. Model

In wildlife studies, it is common to attach radio transmitters to a number of animals and release them. The animals are then located repeatedly by radio telemetry and categorized as being in one of several habitats. It is apparent that repeated locations on the same animal are not independent samples. Researchers attempt to study differences in habitat usage by animals of different ages and sexes.

We therefore consider a sampling scheme consisting of  $J$  subpopulations or strata (defined by age and sex of the animal). From each subpopulation,  $n_j$  animals are sampled from an unknown population of size  $N_j$ . Each animal selected represents a cluster of observations.

Let  $x_{ijk}$  denote the number of observations in the  $i$ th category (habitat) which came from the  $k$ th sampled cluster (animal) of the  $j$ th subpopulation;  $i=2, \dots, I$ ;  $j=1, 2, \dots, J$ ;  $k=1, 2, \dots, n_j$ .  
Let

$$\hat{x}_{jk} = (x_{1jk}, x_{2jk}, \dots, x_{Ijk})'$$

be the observed vector of frequencies for the  $k$ th sampled cluster in the  $j$ th subpopulation. Assume that  $x_{ijk}$  is distributed as a multinomial distribution with parameters  $x_{+jk}$  and

$$P_{jk} = (P_{1jk}, P_{2jk}, \dots, P_{Ijk})'$$

Define the total sample size on the  $k$ th cluster of the  $j$ th subpopulation as

$$x_{+jk} = \sum_{i=1}^I x_{ijk}, \quad (2.1)$$

and let the fixed total sample size,  $x_{+j}$  for the  $j$ th subpopulation be

$$x_{+j} = \sum_{k=1}^{n_j} x_{+jk}. \quad (2.2)$$

Note that  $x_{+jk}$  represents a random sample with replacement from the  $X_{+jk}$  (unknown) observations. Since  $x_{+jk}$  has a multinomial distribution then the density function is

$$f(x_{jk}; P_{jk}, x_{+jk}) = x_{+jk}! \left( \prod_{i=1}^I P_{ijk}^{x_{ijk}} \right)^{-1} \quad (2.3)$$

Define

$$\hat{\lambda}_{jk} = x_{+jk}^{-1} \bar{x}_{jk} \quad (2.4)$$

as the observed vector of proportions for the  $k$ th cluster of the  $j$ th subpopulation, which is an unbiased estimator of the true proportion of times,  $\lambda_{jk}$ , that the  $k$ th cluster of the  $j$ th subpopulation is seen over the  $I$  categories. Let

$$\pi_j = (\pi_{1j}, \pi_{2j}, \dots, \pi_{Ij})'$$

be the true vector of proportions for the  $j$ th subpopulation. These data can be cross classified into a two way contingency table of dimension  $(IXJ)$  where rows represent the  $I$  categories and columns represent the  $J$  subpopulations. This sort of arrangement is very familiar with categorical data.

Define

$$\pi_j = \sum_{k=1}^{N_j} X_{+jk}^{-1} X_{+jk} P_{jk} \quad (2.5)$$

The vector  $\pi_j$  is a weighted linear combination of the true proportion vectors for the  $N_j$  clusters within the  $j$ th subpopulation, where  $X_{+jk}$  is the total number of observations in the  $k$ th cluster in the  $j$ th subpopulation and

$$X_{+j} = \sum_{k=1}^{N_j} X_{+jk} \quad (2.6)$$

is the total number of observations in the  $j$ th subpopulation. Let

$$x_{ij} = \sum_{k=1}^{n_j} x_{ijk} \quad (2.7)$$

be the total sample size for the  $j$ th subpopulation in category  $i$ , then the expected value (denoted by  $E$ ) of the vector  $\bar{x}_j = (x_{1j}, x_{2j}, \dots, x_{Ij})'$  is

$$E(\bar{x}_j) = E \left\{ \sum_{k=1}^{n_j} E(x_{jk} | n_j) \right\}$$

$$= E \left\{ \sum_{k=1}^{n_j} x_{+jk} \frac{P_{jk}}{\lambda_{jk}} \right\}$$

$$E(\bar{x}_j) = \sum_{k=1}^{N_j} n_j X_{+j}^{-1} X_{+jk} x_{+jk} P_{jk} \quad (2.8)$$

Define the  $I$  dimensional observed vector of proportions for the  $j$ th subpopulation as

$$\hat{\lambda}_j = x_{+j}^{-1} \bar{x}_j, \quad (2.9)$$

then under the usual contingency table assumption of fixed marginals,

$$\begin{aligned} E(\hat{\lambda}_j) &= \sum_{k=1}^{N_j} n_j X_{+j}^{-1} X_{+jk} x_{+jk}^{-1} x_{+jk} P_{jk} \\ &= \sum_{k=1}^{N_j} n_j \alpha_{jk} \hat{\alpha}_{jk} P_{jk} \end{aligned} \quad (2.10)$$

where

$$\alpha_{jk} = X_{+j}^{-1} X_{+jk}, \quad (2.11)$$

is an unknown constant and

$$\hat{\alpha}_{jk} = x_{+j}^{-1} x_{+jk}. \quad (2.12)$$

is considered known.

So  $\hat{\lambda}_j$  overestimates the true proportions, since the sum of the weights in (2.10) is greater than one. However, if all the clusters of the same subpopulation are of equal sample size then  $\pi_j$  is an unbiased estimator of  $\pi_j$ . Such an equality condition is rather difficult to satisfy in practice, especially in the study of Wild life. Theoretically, the vector  $\hat{\lambda}_j$  is a type of combined ratio estimator and is expected to be biased (Cochran 1977).

The covariance matrix for the vector  $\hat{\lambda}_j$  conditional on the sample size chosen is

$$\begin{aligned} \text{Var}(\hat{\lambda}_j) &= x_{+j}^{-2} \left\{ \text{Var} \left( \sum_{k=1}^{n_j} x_{+jk} P_{jk} \right) + \right. \\ &\quad \left. E \left( \sum_{k=1}^{n_j} x_{+jk} (\Delta_{P_{jk}} - P_{jk} P_{jk}') \right) \right\} \\ &= x_{+j}^{-2} \left\{ \sum_{k=1}^{N_j} [x_{+jk}^2 n_j \alpha_{jk} (1 - \alpha_{jk}) P_{jk} \right. \\ &\quad \left. P_{jk}' + x_{+jk} n_j \alpha_{jk} (\Delta_{P_{jk}} - P_{jk} P_{jk}') \right\} \\ &= n_j x_{+j}^{-1} \left\{ \sum_{k=1}^{N_j} \hat{\alpha}_{jk} \alpha_{jk} [\Delta_{P_{jk}} + \right. \\ &\quad \left. (x_{+jk} - x_{+jk} \alpha_{jk} - 1) P_{jk} P_{jk}'] \right\}. \end{aligned} \quad (2.13)$$

where  $\Delta_{P_{jk}}$  is a diagonal matrix with elements  $P_{jk}$ .

A consistent estimator of  $\text{Var}(\hat{\lambda}_j)$  is given by  $v(\hat{\lambda}_j)$  where

$$\begin{aligned} v(\hat{\lambda}_j) &= n_j x_{+j}^{-1} \left\{ \sum_{k=1}^{N_j} \hat{\alpha}_{jk} [\Delta_{P_{jk}} + \right. \\ &\quad \left. (x_{+jk} - x_{+jk} \hat{\alpha}_{jk} - 1) \hat{P}_{jk} \hat{P}_{jk}'] \right\}. \end{aligned} \quad (2.14)$$

and defined as  $\hat{B}_j$ , where  $\hat{P}_{jk}$  is an unbiased estimator of  $P_{jk}$ .

### 3. Contingency Table

Let  $(x_{ij})$  denote the contingency table formed using the frequencies  $x_{ij}$  ( $i=1, 2, \dots, I, j=1, 2, \dots, J$ ). The estimator  $\hat{\pi}_{ij}$  which is a multiple of  $x_{ij}$  is not an unbiased estimator of  $\pi_{ij}$  (2.5) unless the  $x_{+jk}$ 's are equal for all  $k$ . The estimator,  $\hat{\pi}_{ij}$  over-estimates  $\pi_{ij}$  except when the sample sizes of the chosen clusters of the  $j$ th subpopulation are all equal.

Consider another two way table denoted by  $(y_{ij})$  formed using the values  $y_{ij}$  ( $i=1, 2, \dots, I; j=1, 2, \dots, J$ ); where

$$y_{ij} = x_{+j} n_j^{-1} \sum_{k=1}^{n_j} x_{ijk} x_{+jk}^{-1} \quad (3.1)$$

$$= x_{+j} n_j^{-1} \sum_{k=1}^{n_j} \hat{p}_{ijk}$$

Define the estimator

$$\hat{\pi}_{ij} = y_{+j}^{-1} y_{ij} \quad (3.2)$$

where  $y_{+j} = (y_{1j}, y_{2j}, \dots, y_{Ij})'$ , and

$$y_{+j} = \sum_{i=1}^I y_{ij}$$

$$= x_{+j} \quad (3.3)$$

The estimator  $\hat{\pi}_{ij}$  is a self weighting estimator, a desirable property in sampling. The estimator  $\hat{\pi}_{ij}$  is a type of separate ratio estimator and is expected to perform well when the relation between  $x_{ijk}$  and  $x_{+jk}$  is constant for a given  $i$  and  $j$ . The expected value of  $\hat{\pi}_{ij}$  conditional on the sample size chosen is

$$E(\hat{\pi}_{ij}) = n_j^{-1} E\left\{ \sum_{k=1}^{n_j} \hat{p}_{ijk} \right\}$$

$$= \pi_{ij} \quad (3.4)$$

Hence  $\hat{\pi}_{ij}$  is an unbiased estimator of  $\pi_{ij}$  regardless of the differences among the clusters' sample sizes. The covariance matrix for the vector,  $\hat{\pi}_{ij}$  conditional on the sample size chosen is

$$\text{Var}(\hat{\pi}_{ij}) = n_j^{-2} \left\{ \text{Var} \sum_{k=1}^{n_j} \hat{p}_{ijk} + E\left( \sum_{k=1}^{n_j} x_{+jk}^{-1} (\Delta_{\hat{p}_{ijk}} - \hat{p}_{ijk} \hat{p}'_{ijk}) \right) \right\}$$

$$= n_j^{-1} \left\{ \sum_{k=1}^{n_j} \alpha_{jk} (1 - \alpha_{jk}) \hat{p}_{ijk} \hat{p}'_{ijk} + x_{+jk}^{-1} \alpha_{jk} (\Delta_{\hat{p}_{ijk}} - \hat{p}_{ijk} \hat{p}'_{ijk}) \right\}$$

$$= n_j^{-1} \sum_{k=1}^{n_j} \alpha_{jk} x_{+jk}^{-1} \{ \Delta_{\hat{p}_{ijk}} + (x_{+jk} - x_{+jk} \alpha_{jk}^{-1}) \hat{p}_{ijk} \hat{p}'_{ijk} \}$$

$$= C_j \quad (3.5)$$

A consistent estimator of  $\text{Var}(\hat{\pi}_{ij})$  is given by  $v(\hat{\pi}_{ij})$  where

$$v(\hat{\pi}_{ij}) = n_j^{-1} \sum_{k=1}^{n_j} x_{+jk}^{-1} \{ \Delta_{\hat{p}_{ijk}} + (x_{+jk} - x_{+jk} \hat{\alpha}_{ijk}^{-1}) \hat{p}_{ijk} \hat{p}'_{ijk} \} \quad (3.6)$$

and defined as  $C_j$ .

The difference between the variances of  $\hat{\pi}_{ij}$  and  $\hat{\pi}_{ij}^*$  is

$$\text{Var}(\hat{\pi}_{ij}) - \text{Var}(\hat{\pi}_{ij}^*) = \sum_{k=1}^{n_j} \alpha_{jk} (n_j x_{+j}^{-1} \hat{\alpha}_{ijk} - n_j^{-1} x_{+jk}^{-1}) v_{jk}, \quad (3.7)$$

where

$$v_{jk} = \Delta_{\hat{p}_{ijk}} + (x_{+jk} - x_{+jk} \alpha_{jk}^{-1}) \hat{p}_{ijk} \hat{p}'_{ijk} \quad (3.8)$$

Hence, the variances are the same whenever  $x_{+jk} = a_j$  for all  $k$ , but it is not necessary that  $a_1 = a_2 = \dots = a_j$ .

Define the coefficient of  $\alpha_{ijk} v_{jk}$  in the difference of variances in (3.8) as

$$R_{jk} = n_j x_{+j}^{-1} \hat{\alpha}_{ijk} - n_j^{-1} x_{+jk}^{-1}$$

$$= n_j x_{+jk} x_{+j}^{-2} - n_j^{-1} x_{+jk}^{-1} \quad (3.9)$$

Since

$$n_j^2 x_{+jk}^2 - x_{+j}^2 = (n_j x_{+jk} - x_{+j}) (n_j x_{+jk} + x_{+j}),$$

the sign of  $R_{jk}$  is unknown.  $R_{jk}$  may be negative, positive or zero for any  $j$ th cluster of the  $j$ th subpopulation. Thus a clear comparison between the variances is not possible. However, Cochran (1977) shows that unless those clusters are really alike the use of  $\hat{\pi}_{ij}$  is likely to be more precise if the sample size in each cluster is large enough to allow estimation of the variance.

Using the table  $(y_{ij})$  with its non integer cell values results in estimators that are unbiased but not necessarily having smaller or larger variances than estimators obtained using table  $(x_{ij})$ . The column totals of table  $(y_{ij})$  are the same as the column totals of table  $(x_{ij})$ . The estimator based on table  $(y_{ij})$  requires a knowledge of the separate total  $x_{+jk}$ , whereas the estimator based on table  $(x_{ij})$  do not require a knowledge of those totals.

### 4. Test of Homogeneity

Consider testing the hypothesis

$$H_0: \pi_j = \pi_0 \quad j=1, 2, \dots, J; \quad (4.1)$$

where  $\pi_0$  is a known vector and  $\pi_j$  is the true vector of proportions for the  $j$ th subpopulation. Then, a Wald type test can be formed

with the biased estimator  $\hat{\pi}_j - \pi_0$  and the covariance  $B_j$ , where  $B_j$  is a consistent estimator of  $B_j$  (2.13). Such a test statistic is given by

$$X_{1H}^2 = \sum_{j=1}^J (\hat{\pi}_j - \pi_0)' B_j^{-1} (\hat{\pi}_j - \pi_0), \quad (4.2)$$

and is asymptotically distributed as a chi-square random variable with  $J(I-1)$  degrees of freedom under  $H_0$ , (Stroud 1971). Similarly, another statistic can be computed using the unbiased estimator  $\tilde{\pi}_j - \pi_0$  and the consistent estimator  $C_j$ . Such a test statistic for testing (4.1) is given by

$$X_{2H}^2 = \sum_{j=1}^J (\tilde{\pi}_j - \pi_0)' \hat{C}_j^{-1} (\tilde{\pi}_j - \pi_0), \quad (4.3)$$

where  $\hat{C}_j$  is a consistent estimator of  $C_j$  (3.6).  $X_{2H}^2$  is also asymptotically distributed as a chi-square random variable with  $J(I-1)$  degrees of freedom, (Stroud 1971).

Consider the contingency table formed based on  $(x_{ij})$  and assuming that the frequencies obtained for the  $j$ th subpopulation follows a multinomial distribution then a test statistic for testing the hypothesis in (4.1) is

$$X_{1Ht}^2 = \sum_{j=1}^J x_{+j} \sum_{i=1}^I (\hat{\pi}_{ij} - \pi_{i0})^2 \pi_{i0}^{-1}. \quad (4.4)$$

$X_{1Ht}^2$  is the usual Pearson Statistic for the test of homogeneity, and is used as an approximation to  $X_{2H}^2$ . Similarly for table  $(y_{ij})$  with the same multinomial assumption, and the use of the estimated vector of proportions  $\tilde{\pi}_j$ , we obtain the test statistic

$$X_{2Ht}^2 = \sum_{j=1}^J x_{+j} \sum_{i=1}^I (\tilde{\pi}_{ij} - \pi_{i0})^2 \pi_{i0}^{-1}, \quad (4.5)$$

as an approximation to  $X_{2H}^2$ . It was shown by Rao and Scott (1981) that  $X_{1Ht}^2$  is a conservative test, for testing the hypothesis in (4.1). In practice, the data are usually available in the form of table  $(x_{ij})$ , so  $H_{1Ht}^2$  is easily calculated. The statistics  $X_{1Ht}^2$  and  $X_{2Ht}^2$  are obtained from the data in the summarized tables  $(x_{ij})$  and  $(y_{ij})$  respectively. They do not require information on each cluster. However, the statistics  $X_{1H}^2$  and  $X_{2H}^2$  cannot be computed from the summarized data given in tables  $(x_{ij})$  and  $(y_{ij})$ . These statistics require information on each cluster.

## 5. Test of Independence

Consider testing the hypothesis

$$H_0: \pi_j = \pi_0 \quad j = 1, 2, \dots, J; \quad (5.1)$$

where  $\pi_0$  is an unknown vector and  $\pi_j$  is the true vector of proportions for the  $j$ th

subpopulation. The unknown vector  $\pi_0$  can be estimated by a weighted linear combination of the  $J$  estimated vectors  $\pi_j$ ,  $j = 1, 2, \dots, J$ ; Thus

$$\hat{\pi}_0 = \sum_{j=1}^J \alpha_j \hat{\pi}_j \quad (5.2)$$

for some known weights,  $\alpha_j$ ,  $j = 1, 2, \dots, J$ ; and an estimator  $\tilde{\pi}_j$  given in (2.9). Similarly, one can define

$$\tilde{\pi}_0 = \sum_{j=1}^J \alpha_j \tilde{\pi}_j \quad (5.3)$$

where  $\tilde{\pi}_j$  is an unbiased estimator of  $\pi_j$  as given in (3.2).

The estimator  $\tilde{\pi}_j - \tilde{\pi}_0$  is an unbiased estimator of  $\pi_j - \pi_0$  for fixed  $\alpha_j$ 's. Let  $T_{jj}$  denote the diagonal elements of  $\text{var}(\tilde{\pi}_j - \tilde{\pi}_0)$  and  $T_{jj'}$  denote the off diagonal elements,  $j \neq j', j, j' = 1, 2, \dots, J$ ; Then as shown in Wilson and Koehler (1984) the matrix

$$T_{jj} = C_j - 2\alpha_j C_j + \sum_{\ell=1}^J \alpha_\ell^2 C_\ell \quad (5.4a)$$

and

$$T_{jj'} = -\alpha_j C_j - \alpha_{j'} C_{j'} + \sum_{\ell=1}^J \alpha_\ell^2 C_\ell. \quad (5.4b)$$

Let  $\hat{T}_{jj}$  and  $\hat{T}_{jj'}$  be consistent estimators of  $T_{jj}$  and  $T_{jj'}$ , respectively.  $\hat{T}_{jj}$  and  $\hat{T}_{jj'}$  are obtained by replacing  $C_j$  with  $\hat{C}_j$ . A test statistic for the hypothesis in (5.1) where  $\pi_0$  is an unknown vector, is

$$X_{2I}^2 = (\tilde{\pi}^{(J)} - \tilde{\pi}_0^{(J)})' \hat{M}_{H_0} (\tilde{\pi}^{(J)} - \tilde{\pi}_0^{(J)}) \quad (5.5)$$

where  $\hat{M}_{H_0}$  is a consistent estimator of

$$M_{H_0} = \text{var}(\tilde{\pi}^{(J)} - \tilde{\pi}_0^{(J)}) \quad (5.6)$$

under  $H_0$ .  $\hat{M}_{H_0}^c$  is the Moore Penrose inverse of  $\hat{M}_{H_0}$ , and the vector of vectors,

$$(\tilde{\pi}^{(J)} - \tilde{\pi}_0^{(J)}) = (\tilde{\pi}_1 - \tilde{\pi}_0, \tilde{\pi}_2 - \tilde{\pi}_0, \dots, \tilde{\pi}_J - \tilde{\pi}_0)'. \quad (5.7)$$

The matrix  $\hat{M}_{H_0}$  has diagonal elements  $\hat{T}_{jj}$  and off diagonal elements  $\hat{T}_{jj'}$ . Similarly, a test statistic can be constructed using the biased estimator  $\hat{\pi}_j - \hat{\pi}_0^{(J)}$  and a consistent estimator of the covariance matrix

$$V_{H_0} = \text{var}(\hat{\pi}^{(J)} - \hat{\pi}_0^{(J)}). \quad (5.8)$$

A consistent estimator,  $\hat{V}_H$  is similar to  $\hat{M}_{H_0}$  except that  $\hat{C}_j$  is replaced by  $\hat{B}_j$  (2.13) in (5.3) and (5.4). Thus,

$$X_{1I}^2 = (\hat{\pi}_0^{(J)} - \hat{\pi}_0^{(J)})' \hat{V}_{H_0}^{-1} (\hat{\pi}_0^{(J)} - \hat{\pi}_0^{(J)}) \quad (5.9)$$

is a test statistic for testing  $H_0$  in (5.1). The statistics  $X_{1I}^2$  and  $X_{2I}^2$  are distributed asymptotically as a chi-square random variable with  $(I-1)(J-1)$  degrees of freedom, (Stroud 1971).

Consider using the summarized data in tables  $(x_{ij})$  and  $(y_{ij})$  based on the multinomial assumption. Then the test statistic based on table  $(x_{ij})$  is

$$X_{1It}^2 = \sum_{j=1}^J x_{+j} \sum_{i=1}^I (\hat{\pi}_{ij} - \hat{\pi}_{i0})^2 \hat{\pi}_{i0}^{-1} \quad (5.10)$$

by Rao & Scott (1981)

The  $Z_i$ 's are standard normal variates. The statistic based on table  $(y_{ij})$  is

$$X_{2It}^2 = \sum_{j=1}^J x_{+j} \sum_{i=1}^I (\hat{\pi}_{ij} - \hat{\pi}_{i0})^2 \hat{\pi}_{i0}^{-1} \quad (5.11)$$

The statistic  $X_{2It}^2$  is the usual Pearson Statistic for the test of independence. It is normally used by researchers as an approximate statistic when the covariance matrix cannot be or is too complicated to estimate to construct of the Wald test statistic. The statistic  $X_{2It}^2$  is similar to  $X_{1It}^2$  in structure and is an approximation to the statistic  $X_{2I}^2$ .

In section 6 in the analysis of the Wild life study data the statistics  $X_{2I}^2$ ,  $X_{2It}^2$  and  $X_{1It}^2$  are related by the expression

$$E\{X_{2I}^2\} \leq E\{X_{2It}^2\} \leq E\{X_{1It}^2\}. \quad (5.12)$$

Thus, having the table constructed with  $(y_{ij})$  as the cell values and using the multinomial assumption results in a less conservative test and a better approximation to the Wald test than the usual Pearson statistic, which is obtained from the use of table  $(x_{ij})$ . Hence in the case where the sampling scheme is as described in section 2 and the estimation of the covariance matrix needed in computing the Wald test, is too complicated, a reasonable approximation is obtained by constructing the alternative contingency table  $(y_{ij})$ , and assuming multinomial sampling. These results suggest that one can obtain better results in terms of approximations in making the adjustments to the construction of the table and then using the multinomial assumption. This requires that the researcher be forewarned about the method of summarization.

## 6. Numerical Example

Data from the study of the diel patterns of habitat use by male and female bobcats in southeastern Oklahoma, Rolley and Warde (1985) were analyzed using the test statistics  $X_{2I}^2$ ,  $X_{1It}^2$  and  $X_{2It}^2$  in (5.5), (5.10) and (5.11) respectively. The data are reproduced in Tables 6.1a and 6.1b. There are  $J=2$  subpopulations, male bobcats and female bobcats. For the male subpopulations, there are  $n_1 = 5$  clusters with vector  $\bar{x}_{+1} = (352, 125, 74, 23, 95)'$ . For the female subpopulation there are  $n_2 = 9$  clusters with vector,  $\bar{x}_{+2} = (195, 19, 90, 72, 26, 74, 60, 95, 52)'$ . There are  $I=5$  categories of interest; pine, deciduous, mixed pine, grassfields, and brush. These categories are assumed to be nonoverlapping and well defined.

The contingency tables  $(x_{ij})$  and  $(y_{ij})$  are given in Tables 6.2a and 6.2b respectively. Table 6.2a is the traditional way of constructing a contingency table while Table 6.2b is the alternative technique proposed in this paper and based on a type of separate ratio estimator. Our hypothesis of interest is  $H_0: \pi_j = \pi_0$  ( $j=1, 2$ ) for some unknown  $\pi_0$ . The idea here is to investigate whether or not the male and female bobcats have the same habitat preferences.

TABLE 6.1a

Diel Patterns of Habitat Use by Five Female Bobcats in Southeastern Oklahoma

HABITATS	Bobcats				
	1	2	3	4	5
Pine	227	80	50	9	39
Deciduous	53	10	3	4	0
Mixed Pine	53	30	20	9	11
Grass Fields	8	5	1	0	31
Brush	11	0	0	1	14
Total	352	125	74	23	95

TABLE 6.1b

Diel Patterns of Habitat Use by Nine Male Bobcats in Southeastern Oklahoma

HABITATS	1	2	3	4	5	6	7	8	9
	Pine	145	11	49	38	6	33	46	39
Deciduous	3	1	11	11	13	3	0	5	8
Mixed Pine	26	2	21	15	6	18	9	14	13
Grass Fields	11	4	4	2	1	1	2	30	1
Brush	10	1	5	6	0	19	3	7	11
Total	195	19	90	72	26	74	60	95	52

Under Table 6.2a the estimated vectors are for female bobcats  $\hat{\pi}_1 = (.605, .105, .184, .067, .039)'$  and for male bobcats  $\hat{\pi}_2 = (.565, .081, .181, .082, .091)'$ . The usual Pearson statistic  $X_{1It}^2$  given in (5.10), is 18.289. Under Table 6.2b the estimated vector for female bobcats,  $\hat{\pi}_1 = (.553, .089, .234, .081, .044)'$  and for male bobcats,  $\hat{\pi}_2 = (.513, .121, .189, .084, .093)'$ . The alternative statistic  $X_{2It}^2$  given in (5.11) is 9.171. From Tables 6.1a and 6.1b the statistics  $X_{1I}^2$  in (5.9) and  $X_{2I}^2$  in (5.5)

were calculated. These are Wald statistics. The diagonal elements of the covariance matrix used in calculating  $X_{1I}^2$  are (.01283, .00032, .00177, .00706, .00209, .00863, .00096, .00271, .00233, .00224)'. The diagonal elements of the covariance matrix used in calculating  $X_{2I}^2$  are (.02574, .00069, .00363, .01405, .00425, .01786, .00203, .00573, .00464, .00470)'. The statistic  $X_{2I}^2$  based on the separate type estimator has the value 6.664 and the statistic  $X_{1I}^2$  based on the combined type estimator has the value 8.830. Statistics  $X_{1It}^2$  and  $X_{2It}^2$  are the approximations to  $X_{2I}^2$ . While  $X_{1It}^2$  is an unsuitable approximation the statistic  $X_{2It}^2$  is a reasonable estimator.

When these statistics are considered to be distributed as chi-square random variables with 4 degrees of freedom, we are led to rejecting the null hypothesis at the 5% significant level, if we use  $X_{1It}^2$ , the usual Pearson Statistic. All other statistics considered in this example led to supporting the claim that the bobcats (males and females) have about the same habitat preference in Southeastern Oklahoma.

TABLE 6.2a

A Habitat by Sex Two Way Contingency Table for Bobcats in Southeastern Oklahoma

	Females	Males
Pine	405	386
Deciduous	70	55
Mixed Pine	123	124
Grassfields	45	56
Brush	26	62

TABLE 6.2b

A Habitat by Sex Two Way Alternative Contingency Table for Bobcats in Southeastern Oklahoma

	Females	Males
Pine	369.689	350.227
Deciduous	59.541	82.719
Mixed Pine	156.278	129.011
Grassfields	53.921	57.524
Brush	29.570	63.443

## 7. DISCUSSION

The presence of clustering in the collection of sample data can have a severe effect on certain test statistics obtained from the frequency data in a usual contingency table. Such computed statistics are usually too large in numerical value. A better approximation is to construct the table based on a separate type estimator and then to use the usual techniques of constructing Pearson statistics. This technique has its greatest gain when the clusters differ greatly.

Rao and Scott (1981, 1984), Bedrick (1983), Wilson and Koehler (1984), Brier (1980) and Holt, Scott and Ewings (1980) have considered model that leads to a correction of the usual Pearson Statistics. Their works rely on

summarized data through the usual construction of a contingency table. However, in this paper no correction is considered for the usual Pearson Statistic. Here the changes are suggested prior to the summarized data. There is no need for matrix inversion or the computation of several covariances. Eigen values are not needed. The computer programs necessary are readily available. They are the same as when multinomial sampling is conducted and a Pearson Statistic computed.

## REFERENCES

- Bedrick, E.J. (1983). Adjusted Chi-squares tests for cross-classified tables of survey data. *Biometrika* 70: 591-595.
- Brier, S.S. (1980). Analysis of Contingency tables under cluster Sampling. *Biometrika* 67: 591-596.
- Cohen, J.E. (1976). The distribution of the chi-square statistic under clustered sampling from contingency tables. *J. Amer. Statist. Assoc.* 71: 665-670.
- Cochran, W.G. (1977). *Sampling Techniques*. John Wiley & Sons, 3rd Edition, New York.
- Holt, D., Scott, A.J., and Ewings, P.D. (1980). Chi-squared tests with survey data. *J.R. Statist. Soc., Series A.* 143 (3): 303-320.
- Rao, J.N.K. and Scott, A.J. (1981). The Analysis of Categorical Data from complex surveys. *J. Amer. Statist. Assoc.* 76: 221-30.
- Rao, J.N.K. and Scott, A.J. (1984). On Chi-squared test for Multiway Contingency Tables with all proportions estimated from survey data. *The Annals of Statistics* 1984, Vol. 2, No. 1, 46-60.
- Rolley, R.E. and Warde, W.D. (1985). Bobcat habitat use and activity patterns in Southeastern Oklahoma. *Journal of Wildlife Management* (in press).
- Stroud, T.W. F. (1971). On obtaining large sample tests from asymptotic normal estimations. *Annals of Mathematic Statistics Association*, 72, 881-885.
- Wald, A. (1943). Tests of Statistical hypotheses concerning several parameters when the number of observations is large. *Trans. Amer. Math. Soc.* 54: 426-482.
- Wilson, J.R. (1984). *Statistical Methods for frequency data from complex sampling schemes*. Ph.D. Dissertation, Iowa State University, Ames, Iowa 50011.
- Wilson, J.R. and Koehler, K.J. (1984). Testing equality of vectors of proportions for several cluster samples. *Proceedings of Joint Statistical Association Meetings 1984. Section on Survey Research Methods*, 201-206.