

This research was supported by National Institutes of Health, Institute of General Medical Sciences Grants GM-70004-01 and GM-12868-08.

AN ANALYSIS FOR COMPOUNDED LOGARITHMIC-EXPONENTIAL-
LINEAR FUNCTIONS OF CATEGORICAL DATA

by

Ronald N. Forthofer and Gary G. Koch

Department of Biostatistics
University of North Carolina at Chapel Hill

Institute of Statistics Mimeo Series No. 801

February 1972

AN ANALYSIS FOR COMPOUNDED LOGARITHMIC - EXPONENTIAL -
LINEAR FUNCTIONS OF CATEGORICAL DATA

by

Ronald N. Forthofer
School of Public Health
University of Texas
Houston, Texas 77025, U.S.A.

Gary G. Koch
Department of Biostatistics
University of North Carolina
Chapel Hill, N.C., 27514, U.S.A.

1. INTRODUCTION

One area of application which has become increasingly important to statisticians and other researchers is the analysis of categorical data. Often the principal objective in such investigations is either the testing of appropriate hypotheses or the fitting of simplified models to the multi-dimensional contingency tables which arise when frequency counts are obtained for the respective cross-classifications of specific qualitative variables. Grizzle, Starmer, and Koch [1969] (subsequently abbreviated GSK) have described how linear regression models and weighted least squares can be used for this purpose. The resulting test statistics belong to the class of minimum modified chi-square due to Neyman [1949] which is equivalent to the general quadratic form criteria of Wald [1943]. As such, they have central χ^2 -distributions when the corresponding null hypotheses are true. Two alternative approaches to this methodology are that based on maximum likelihood as formulated by Bishop [1969, 1971] and Goodman [1970, 1971a, 1971b] and that based on minimum discrimination information as formulated by Ku, Varner, and Kullback [1971].

In each of the previously mentioned papers, primary emphasis was given to the formulation of models and the problems of analysis under various conditions of "no interaction" (see Roy and Kastenbaum [1956] or Bhapkar and Koch [1968a, 1968b]) either in the overall contingency table or the marginals thereof. As a result, attention was essentially restricted to either linear or log-linear functions of cell probabilities. However, there exist a number of common practical situations in which more complex functions must be handled. Examples here include:

1. The analysis of patterns of association in square contingency tables as related to functions of diagonal totals
2. The analysis of rank correlation coefficients as discussed by Goodman and Kruskal [1954, 1959, 1963] and Davis and Quade [1968]
3. The analysis of "ridits" as discussed by Bross [1958]
4. The analysis of partial association as discussed by Mantel and Haenszel [1959] and Mantel [1963].

Each of these problems can be formulated in terms of compounded estimators involving logarithmic, exponential, and linear functions. Thus, such analyses can be undertaken by using an appropriate extension of the GSK approach to this broader class of functions. The remainder of this paper will be concerned with describing and illustrating the use of this more general methodology.

2. STATISTICAL THEORY

We shall assume that there are s populations of elements from which independent random samples of fixed sizes n_1, n_2, \dots, n_s respectively are selected. The responses of the n_i elements from the i -th population are classified into r categories with n_{ij} , where $j = 1, 2, \dots, r$ denoting the number of elements classified in the j -th response category for the i -th population.

The vector \tilde{n}_i , where $\tilde{n}_i' = (n_{i1}, n_{i2}, \dots, n_{ir})$, will be assumed to follow the multinomial distribution with parameters n_i and $\tilde{\pi}_i' = (\pi_{i1}, \pi_{i2}, \dots, \pi_{ir})$.

Thus, the relevant product multinomial model is

$$\phi = \prod_{i=1}^s \left\{ \frac{n_i!}{\prod_{j=1}^r n_{ij}!} \right\} \prod_{j=1}^r \pi_{ij}^{n_{ij}} \quad (2.1)$$

Let $\underline{p}_i = (n_i/n_i)$ and let \underline{p} be the compound vector defined by $\underline{p}' = (p'_1, p'_2, \dots, p'_s)$. A consistent estimate for the covariance matrix of \underline{p} is given by the block diagonal matrix $\underline{V}(\underline{p})$ with the matrices $\underline{V}_i(\underline{p}_i) = (D_{\underline{p}_i} - \underline{p}_i \underline{p}'_i)/n_i$ for $i = 1, 2, \dots, s$ representing the main diagonal; here $D_{\underline{p}_i}$ is a diagonal matrix with elements in the vector \underline{p}_i on the main diagonal.

Let $F_1(\underline{p}), F_2(\underline{p}), \dots, F_u(\underline{p})$ be a set of u functions of \underline{p} , each with partial derivatives up to order two with respect to the elements of \underline{p} existing throughout a region containing $\underline{\pi} = \underline{\varepsilon}(\underline{p})$ where $\underline{\pi}' = (\pi'_1, \pi'_2, \dots, \pi'_s)$. If $\underline{F} \equiv \underline{F}(\underline{p})$ is defined by $\underline{F}' = \underline{F}'(\underline{p}) = (F_1(\underline{p}), F_2(\underline{p}), \dots, F_u(\underline{p}))$, then the covariance matrix of \underline{F} can be consistently estimated by $\underline{V}_{\underline{F}} = \underline{H}[\underline{V}(\underline{p})]\underline{H}'$ where $\underline{H} = [d\underline{F}(\underline{x})/d\underline{x} | \underline{x} = \underline{p}]$. In all applications, the functions comprising \underline{F} are chosen so that $\underline{V}_{\underline{F}}$ is asymptotically non-singular.

The function vector \underline{F} is a consistent estimator of $\underline{F}(\underline{\pi})$. Hence, consideration can then be directed at fitting a linear model $\underline{\varepsilon}\{\underline{F}(\underline{p})\} \doteq \underline{F}(\underline{\pi}) = \underline{X}\underline{\beta}$, where \underline{X} is a known $(u \times t)$ coefficient matrix of full rank $t \leq u$ and $\underline{\beta}$ is an unknown $(t \times 1)$ parameter vector. Weighted least squares is applied to determine a BAN estimator \underline{b} for $\underline{\beta}$ as obtained from the expression

$$\underline{b} = (\underline{X}'\underline{V}_{\underline{F}}^{-1}\underline{X})^{-1}\underline{X}'\underline{V}_{\underline{F}}^{-1}\underline{F} \quad (2.2)$$

A consistent estimate for the covariance matrix of \underline{b} is given by

$$\underline{V}_{\underline{b}} = (\underline{X}'\underline{V}_{\underline{F}}^{-1}\underline{X})^{-1} \quad (2.3)$$

A goodness of fit statistic for assessing the extent to which the model characterizes the data is the residual sum of squares

$$X^2 = SS(\underline{\varepsilon}\{\underline{F}\} = \underline{X}\underline{\beta}) = \underline{F}'\underline{V}_{\underline{F}}^{-1}\underline{F} - \underline{b}'(\underline{X}'\underline{V}_{\underline{F}}^{-1}\underline{X})\underline{b} \quad (2.4)$$

which has approximately a chi-square distribution with D.F. = $(u - t)$ in large samples under the hypothesis that the model fits. If the model does adequately

describe the data, tests of linear hypotheses with respect to the parameters comprising β can be undertaken. In particular, for a general hypothesis of the form $H_0: C\beta = 0$ where C is a known $(d \times t)$ matrix of full rank $d \leq t$, a suitable test statistic is

$$X^2 = SS(C\beta = 0) = b' C' [C(X' V_F^{-1} X)^{-1} C']^{-1} C b \quad (2.5)$$

which has approximately a chi-square distribution with D.F. = d in large samples under H_0 .

Although this formulation is quite general, GSK restricted attention to two types of functions. These were linear functions of the type

$$F(p) = Ap \equiv a \quad (2.6)$$

where A is a known $(u \times rs)$ matrix and log-linear functions of the type

$$F(p) = K[\log_e(Ap)] = K[\log_e(a)] = f \quad (2.7)$$

where K is a known $(k \times u)$ matrix, A is as defined in (2.6), and \log_e transforms a vector to the corresponding vector of natural logarithms. For the linear functions in (2.6), the corresponding estimated covariance matrix is

$$V_a = A[V(p)]A' \quad (2.8)$$

while for the log-linear functions in (2.7), the corresponding estimated covariance matrix is

$$V_f = K D_a^{-1} A[V(p)] A' D_a^{-1} K' \quad (2.9)$$

where D_a is a diagonal matrix with elements of the vector a on the main diagonal.

In this paper, two other general classes of functions will be considered.

The first of these are exponential functions of the type

$$F(p) = Q(\exp\{K[\log_e(Ap)]\}) = Q(\exp\{f\}) = g \quad (2.10)$$

where Q is a known $(q \times k)$ matrix, A is as defined in (2.6), K is as defined in (2.7), and \exp transforms a vector to the corresponding vector of exponential functions

(i.e., of anti-logarithms). A consistent estimate for the covariance matrix of \underline{g} is given in (2.11) where $\underline{y} = \exp(\underline{f})$.

$$\underline{V}_{\underline{g}} = \underline{Q} \underline{D} \underline{K} \underline{D}^{-1} \underline{A} [\underline{V}(\underline{p})] \underline{A}' \underline{D}^{-1} \underline{K}' \underline{D} \underline{Q}' \quad (2.11)$$

The other class of functions are compounded logarithmic functions of the type

$$\underline{F}(\underline{p}) = \underline{L} \{ \log_e [\underline{Q} (\exp \{ \underline{K} [\log_e (\underline{A} \underline{p})] \})] \} = \underline{L} \{ \log_e (\underline{g}) \} = \underline{h} \quad (2.12)$$

where \underline{L} is a known $(l \times q)$ matrix, \underline{A} is defined in (2.6), \underline{K} is defined in (2.7), and \underline{Q} is defined in (2.10). A consistent estimate for the covariance matrix of \underline{h} is

$$\underline{V}_{\underline{h}} = \underline{L} \underline{D}^{-1} \underline{Q} \underline{D} \underline{K} \underline{D}^{-1} \underline{A} [\underline{V}(\underline{p})] \underline{A}' \underline{D}^{-1} \underline{K}' \underline{D} \underline{Q}' \underline{D}^{-1} \underline{L}' \quad (2.13)$$

From the results in (2.10) and (2.11) together with those in (2.12) and (2.13), it is apparent that logarithmic, exponential, and linear functions can be compounded to a further extent if necessary. Thus, a considerably broader class of functions can be readily analyzed by the use of linear models than was originally demonstrated by GSK.

Finally, this extended methodology can be undertaken by using the same computer program cited in GSK. The only refinements required are suitable options for additional matrix operations. Also, in some cases, if certain n_{ij} are zero, it may be necessary to replace them by $(1/rn_i)$ in order to prevent $\underline{V}_{\underline{F}}$ from being singular (see GSK and Berkson [1955] for more details in this respect).

3. EXAMPLES

3.1. An analysis of patterns of association related to diagonals

Let us reconsider the first example discussed by GSK which pertained to the unaided distance vision of women aged 30-39. However, here we shall

analyze the contingency table shown in Table 1 in which the frequencies have been reduced to approximately 10% of their original magnitude.

Table 1
Unaided Distance Vision; 747 Women 30-39

<u>Right Eye</u>	<u>Left Eye</u>				Total
	Highest Grade (1)	Second Grade (2)	Third Grade (3)	Lowest Grade (4)	
Highest Grade (1)	152	27	12	7	198
Second Grade (2)	23	151	43	8	225
Third Grade (3)	12	36	177	20	245
Lowest Grade (4)	4	8	18	49	79
Total	191	222	250	84	747

One question of interest with respect to these data was whether there was a difference between right eye vision and left eye vision which corresponds to a hypothesis of marginal homogeneity. Various test statistics for this condition have been discussed by several authors. On the other hand, there is also some interest in the association between right eye vision and left eye vision. In particular, if there is a strong relationship between vision in the two eyes, then the frequencies along the main diagonal (i.e., cells 11, 22, 33, 44) should be substantially inflated over what would be expected if right eye vision and left eye vision were independent. Moreover, this inflation might also smoothly decrease in a symmetric manner to depressions with respect to other such diagonals which were successively farther above and/or below the main diagonal. Although this formulation of a pattern of association appears complex, the corresponding analysis can be undertaken in a straightforward manner by using the

methodology outlined in Section 2.

Since the data come from a single multinomial population, the vector \underline{p} is defined by expression (3.1);

$$\underline{p}' = [152 \ 27 \ 12 \ 7 \ 23 \ 151 \ 43 \ 8 \ 12 \ 36 \ 177 \ 20 \ 4 \ 8 \ 18 \ 49]/747; \quad (3.1)$$

the corresponding estimated covariance matrix is given by $\underline{V}(\underline{p}) = (\underline{D}_{\underline{p}} - \underline{p}\underline{p}')/747$. The matrix \underline{A} in (3.2) is applied to \underline{p} to form diagonal proportions together with suitable marginal proportions which are used to form estimators for the expected values of the diagonals under the hypothesis of independence. Appropriate pairs

$$\underline{A} = \begin{matrix} \underline{A}_1 \\ \underline{A}_2 \end{matrix} \begin{matrix} 5 \times 16 \\ 8 \times 16 \end{matrix} \quad \text{where} \quad \underline{A}_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \quad (3.2)$$

$$\underline{A}_2 = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

of these marginal proportions are in effect multiplied together on the log-scale by the \underline{K} -matrix in (3.3) where \underline{J}_4 is a (4×1) vector of 1's

$$\underline{K} = \begin{matrix} 21 \times 13 \end{matrix} = \begin{bmatrix} \underline{I}_5 & \underline{0}_5 & \underline{0}_5 & \underline{0}_5 & \underline{0}_5 & \underline{0}_{5,4} \\ \underline{0}_{4,5} & \underline{J}_4 & \underline{0}_4 & \underline{0}_4 & \underline{0}_4 & \underline{I}_4 \\ \underline{0}_{4,5} & \underline{0}_4 & \underline{J}_4 & \underline{0}_4 & \underline{0}_4 & \underline{I}_4 \\ \underline{0}_{4,5} & \underline{0}_4 & \underline{0}_4 & \underline{J}_4 & \underline{0}_4 & \underline{I}_4 \\ \underline{0}_{4,5} & \underline{0}_4 & \underline{0}_4 & \underline{0}_4 & \underline{J}_4 & \underline{I}_4 \end{bmatrix} \quad (3.3)$$

0_5 is a (5×1) vector of 0's, $0_{5,4}$ is a (5×4) matrix of 0's, I_5 is the (5×5) identity matrix, etc. Corresponding values for observed and expected proportions for the diagonals are obtained with the Q matrix in (3.4) with A_1 as in (3.2).

Finally,

$$Q = \begin{matrix} 10^{21} \times 21 \\ \begin{bmatrix} I_5 & 0_{5,16} \\ 0_{5,5} & A_1 \end{bmatrix} \end{matrix} \quad (3.4)$$

the appropriate ratios of these estimates on the log-scale result on using

$$\begin{matrix} L \\ 5 \times 10 \end{matrix} = \begin{bmatrix} I_5 & -I_5 \end{bmatrix} \quad (3.5)$$

For the data in Table 1, the function vector F obtained from expression (2.12) with A , K , Q , and L defined by (3.2), (3.3), (3.4), and (3.5) respectively is

$$F = \begin{bmatrix} \log_e \{(p_{11}+p_{22}+p_{33}+p_{44})/(p_{10}p_{01}+p_{20}p_{02}+p_{30}p_{03}+p_{40}p_{04})\} \\ \log_e \{(p_{12}+p_{23}+p_{34})/(p_{10}p_{02}+p_{20}p_{03}+(p_{30}p_{04}))\} \\ \log_e \{(p_{13}+p_{24})/(p_{10}p_{03}+p_{20}p_{04})\} \\ \log_e \{(p_{21}+p_{32}+p_{43})/(p_{20}p_{01}+p_{30}p_{02}+p_{40}p_{03})\} \\ \log_e \{(p_{31}+p_{42})/(p_{30}p_{01}+p_{40}p_{02})\} \end{bmatrix} = \begin{bmatrix} 0.931 \\ -0.576 \\ -1.523 \\ -0.712 \\ -1.462 \end{bmatrix} \quad (3.6)$$

given in (3.6) where $p_{jj'}$ is the element in p corresponding to the j -th category of right eye vision and the j' -th category of left eye vision, $p_{j.} = \sum_{j'=1}^4 p_{jj'}$, and $p_{.j'} = \sum_{j=1}^4 p_{jj'}$. On applying (2.13), an estimated covariance matrix for F is the matrix V_F in (3.7).

$$V_F = \begin{bmatrix} 0.74 & -1.20 & -1.26 & -1.27 & -1.38 \\ & 8.56 & -3.64 & -0.17 & 1.03 \\ & & 43.42 & 1.10 & 3.69 \\ & & & 10.28 & -3.87 \\ & & & & 43.99 \end{bmatrix} \times 10^{-3} \quad (3.7)$$

Finally, it is appropriate to note that diagonal functions corresponding to the upper right hand corner cell (i.e., $\log_e \{p_{14}/p_{10}p_{04}\}$) and the lower left hand corner cell (i.e., $\log_e \{p_{41}/p_{40}p_{01}\}$) have been excluded here to prevent possible singularities in \tilde{V}_F .

If right eye vision and left eye vision were independent, then the observed and expected values for the diagonal proportions would be approximately equal and the elements of \tilde{F} would each be approximately zero. Hence, a statistical test for the hypothesis of independence can be obtained by fitting a suitable linear model which explains essentially all the variation among the components of \tilde{F} and then testing whether the parameters of that model all equal zero. For \tilde{F} in (3.6), one model of interest is

$$\tilde{\varepsilon}\{\tilde{F}\} = \tilde{X}\tilde{\beta} \quad (3.8)$$

1	0	0	β_0 β_1 β_2	(3.8)
1	1	1		
1	2	4		
1	1	1		
1	2	4		

where β_0 is an overall mean while β_1 and β_2 reflect linear and quadratic components of the variation of the elements in \tilde{F} with respect to the absolute value of the difference between the grades of vision for the right and left eye where the four grades have been scaled as "1," "2," "3," and "4" as shown in Table 1. For the data in this example, the goodness of fit statistic in (2.4) corresponding to the model (3.8) was $X^2 = 0.95$ with D.F. = 2 which is non-significant ($\alpha = .25$). Thus, this model suitably characterizes the vector \tilde{F} .

The estimate \tilde{b} of $\tilde{\beta}$ for the model (3.8) and its estimated covariance matrix \tilde{V}_b , which have been obtained by using (2.2) and (2.3) are shown in (3.9). Since the hypothesis of independence corresponds to $H_0: \beta_0 = \beta_1 = \beta_2 = 0$, an appropriate

$$\tilde{b} = \begin{bmatrix} 0.930 \\ -1.924 \\ 0.356 \end{bmatrix} \quad \tilde{V}_{\tilde{b}} = \begin{bmatrix} 0.73 & -2.90 & .94 \\ & 34.00 & 19.42 \\ & & 12.62 \end{bmatrix} \times 10^{-3} \quad (3.9)$$

test statistic is (2.5) with $\tilde{C} = \tilde{I}_3$ in which case $X^2 = 1657.22$ with D.F. = 3. Thus, the association between right eye vision and left eye vision is statistically significant ($\alpha = .01$). Moreover, for the hypotheses $H_0: \beta_2 = 0$ and $H_0: \beta_1 = 0$, the results are $X^2 = 10.04$ with D.F. = 1 and $X^2 = 108.83$ with D.F. = 1, both of which are significant ($\alpha = .01$). Hence, the model (3.8) together with its estimated parameters in (3.9) represent a concise and interesting characterization of the pattern of association between right eye vision and left eye vision as reflected by these data.

3.2. An analysis of rank correlation

In Table 2, data are shown for the severity of the 'dumping syndrome,' an undesirable sequela of surgery for duodenal ulcer where the four surgical operations are labelled as

- A = Drainage and vagotomy
- B = 25% resection (antrectomy) and vagotomy
- C = 50% resection (hemigastrectomy) and vagotomy
- D = 75% resection.

These data were previously considered by GSK in terms of an analysis of variance of the mean scores corresponding to the respective surgical procedure \times hospital combinations with none, slight, and moderate being scaled as 1, 2, and 3. Alternatively, the association between the severity of the dumping syndrome and the surgical procedure can be investigated by using rank correlation coefficients described by Goodman and Kruskal [1954, 1959, 1963] and Davis and Quade [1968] since both of these variables are ordinal. Such analyses can be formulated in

terms of the methodology of Section 2.

Table 2

Hospital	Operation A			Operation B			Operation C			Operation D		
	None	Slight	Mod.	None	Slight	Mod.	None	Slight	Mod.	None	Slight	Mod.
I	23	7	2	23	10	5	20	13	5	24	10	6
II	18	6	1	18	6	2	13	13	2	9	15	2
III	8	6	3	12	4	4	11	6	2	7	7	4
IV	12	9	1	15	3	2	14	8	3	13	6	4
Total	61	28	7	68	23	13	58	40	12	53	38	16

First of all, the data for each hospital will be assumed to have arisen from separate and independent multinomial populations. Let $p_i = (n_{ij}/n_i)$ where n_i is the vector of frequencies in the i -th row of Table 2 and n_i is their sum, $i = I, II, III, IV$. The A matrix in (3.10) produces all cell probabilities together with sums of appropriate concordant and discordant combinations for each of the separate hospitals. Then, the corresponding probabilities of the various types of concordant

$$A = \begin{matrix} \begin{matrix} I_{12} \\ \sim 1 \\ A_1 \end{matrix} \\ \begin{matrix} I_{12} \\ \sim 1 \\ A_1 \end{matrix} \\ \begin{matrix} I_{12} \\ \sim 1 \\ A_1 \end{matrix} \\ \begin{matrix} I_{12} \\ \sim 1 \\ A_1 \end{matrix} \end{matrix} \quad \text{where} \quad \begin{matrix} A_1 = \\ 12 \times 12 \end{matrix} \begin{matrix} 000 & 011 & 011 & 011 \\ 000 & 001 & 001 & 001 \\ 000 & 000 & 011 & 011 \\ 000 & 000 & 001 & 001 \\ 000 & 000 & 000 & 011 \\ 000 & 000 & 000 & 001 \\ 000 & 100 & 100 & 100 \\ 000 & 110 & 110 & 110 \\ 000 & 000 & 100 & 100 \\ 000 & 000 & 110 & 110 \\ 000 & 000 & 000 & 100 \\ 000 & 000 & 000 & 110 \end{matrix} \quad (3.10)$$

and discordant pairs are obtained on the log-scale with the \tilde{K} -matrix in (3.11).

$$\tilde{K} = \begin{matrix} 48 \times 96 \\ \left[\begin{array}{cccc} \tilde{K}_1 & I_{12} & & \\ & \tilde{K}_1 & I_{12} & \\ & & \tilde{K}_1 & I_{12} \\ & & & \tilde{K}_1 & I_{12} \end{array} \right] \end{matrix} \quad \text{where} \quad \tilde{K}_1 = \begin{matrix} 12 \times 12 \\ \left[\begin{array}{cccc} 100 & 000 & 000 & 000 \\ 010 & 000 & 000 & 000 \\ 000 & 100 & 000 & 000 \\ 000 & 010 & 000 & 000 \\ 000 & 000 & 100 & 000 \\ 000 & 000 & 010 & 000 \\ 010 & 000 & 000 & 000 \\ 001 & 000 & 000 & 000 \\ 000 & 010 & 000 & 000 \\ 000 & 001 & 000 & 000 \\ 000 & 000 & 010 & 000 \\ 000 & 000 & 001 & 000 \end{array} \right] \end{matrix} \quad (3.11)$$

These are combined into overall estimates for the probabilities of concordance and discordance by the use of the \tilde{Q} -matrix in (3.12)

$$\tilde{Q} = \begin{matrix} 8 \times 48 \\ \left[\begin{array}{cccc} \tilde{Q}_1 & & & \\ & \tilde{Q}_1 & & \\ & & \tilde{Q}_1 & \\ & & & \tilde{Q}_1 \end{array} \right] \end{matrix} \quad \text{where} \quad \tilde{Q}_1 = \begin{matrix} 2 \times 12 \\ \left[\begin{array}{cccccccc} 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 \end{array} \right] \end{matrix} \quad (3.12)$$

Finally, the ratio of the probabilities of concordance and discordance for each of the four hospitals are obtained with the \tilde{L} -matrix in (3.13). This final transformation produces indices of association which are analogous to the log cross-product ratios in (2x2) tables; also it will be shown later that these indices can be readily transformed into rank correlation coefficients.

For the data in Table 2, the function vector \tilde{F} obtained from expression (2.12)

$$\tilde{L} = \begin{matrix} 4 \times 8 \\ \left[\begin{array}{cccccc} 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{array} \right] \end{matrix} \quad (3.13)$$

with \tilde{A} , \tilde{K} , \tilde{Q} and \tilde{L} defined by (3.10), (3.11), (3.12), and (3.13) respectively is given in (3.14) together with the corresponding estimated covariance matrix $\tilde{V}_{\tilde{F}}$.

$$\tilde{F} = \begin{bmatrix} .272 \\ .813 \\ .137 \\ .169 \end{bmatrix} \quad \tilde{V}_{\tilde{F}} = \begin{bmatrix} 5.0464 & & & \\ & 7.9299 & & \\ & & 8.9247 & \\ & & & 8.4030 \end{bmatrix} \times 10^{-2} \quad (3.14)$$

Once the vector \tilde{F} has been obtained, a question arises as to the extent to which the rank correlation coefficients comprising \tilde{F} do not depend on the hospital. This can be investigated by fitting the linear model in (3.15) where β_0 may

$$\tilde{\epsilon}\{\tilde{F}\} = \tilde{X}\beta = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \beta_0 \quad (3.15)$$

be interpreted as a measure of the partial association between the severity of the dumping syndrome and the extent of the operation. The goodness of fit statistic in (2.4) corresponding to the model (3.15) is $X^2 = 3.72$ with D.F. = 3 which is non-significant ($\alpha = .25$). Hence, the model (3.15) adequately characterizes the vector \tilde{F} which may be interpreted as meaning that there is no second order interaction between hospitals and surgical operations with respect to their effect on the severity of the dumping syndrome. This result is consistent with that obtained by GSK using another method of analysis. Given this result, the estimate b of β_0 in (3.15) and its estimated variance are

$$b = 0.346 \quad V_b = 0.018008. \quad (3.16)$$

For the hypothesis $H_0: \beta_0 = 0$, the test statistic in (2.5) is $X^2 = 6.63$ with D.F. = 1. Thus, this measure of partial association is statistically significant ($\alpha = .01$) which implies that there is a definite relationship between the severity of the dumping syndrome and the surgical operation. This result is also

consistent with that previously obtained by GSK.

Throughout this discussion, the indices \tilde{F} have been considered. However, using additional operations analogous to (2.10) and (2.12), the actual rank correlation coefficients and their estimated variances can be obtained as

$$\tilde{G} = (\tilde{D}_{\tilde{g}} + \tilde{I}_4)^{-1} (\tilde{g} - \tilde{J}_4) \quad (3.17)$$

$$\tilde{V}_{\tilde{G}} = 4(\tilde{D}_{\tilde{g}} + \tilde{I}_4)^{-2} \tilde{D}_{\tilde{g}} \tilde{V}_{\tilde{F}} \tilde{D}_{\tilde{g}} (\tilde{D}_{\tilde{g}} + \tilde{I}_4)^{-2} \quad (3.18)$$

where $\tilde{g} = \exp(\tilde{F})$, $\tilde{D}_{\tilde{g}}$ is a diagonal matrix with elements of \tilde{g} on the diagonal, \tilde{J}_4 is a 4×1 vector of 1's, and \tilde{I}_4 is the 4×4 identity matrix. The resulting rank correlation coefficients \tilde{G} and $\tilde{V}_{\tilde{G}}$ are shown in (3.19). In addition, the partial rank correlation coefficient and its estimated variance can be determined

$$\tilde{G} = \begin{bmatrix} 0.135 \\ 0.385 \\ 0.068 \\ 0.084 \end{bmatrix} \quad \tilde{V}_{\tilde{G}} = \begin{bmatrix} 1.22 & & & \\ & 1.44 & & \\ & & 2.22 & \\ & & & 2.07 \end{bmatrix} \times 10^{-2} \quad (3.19)$$

by similarly treating b and V_b in (3.16) to obtain

$$\bar{G} = (e^b - 1)/(e^b + 1) = 0.171 \quad V_{\bar{G}} = \frac{1}{4}(1 - \bar{G}^2)^2 V_b = 0.004242 \quad (3.20)$$

A more direct approach to obtain \tilde{G} and $\tilde{V}_{\tilde{G}}$ is to use a different Q_1 matrix in (3.12). In particular, if

$$Q_1 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad (3.21)$$

and \tilde{L} is the same as in (3.13), then \tilde{G} in (3.19) can be obtained from \tilde{F} based on (3.21) by using $\tilde{G} = \{2[\exp(\tilde{F})] - \tilde{J}_4\}$. In this event, if $\tilde{g} = 2[\exp(\tilde{F})]$, then

$\tilde{V}_{\tilde{G}} = \tilde{D}_{\tilde{g}} \tilde{V}_{\tilde{F}} \tilde{D}_{\tilde{g}}$. Finally, both of these methods for the computation of rank correlation

coefficients and their variance yield essentially the same results as would be obtained by using the approach of either Goodman and Kruskal [1954,1959,1963] or Davis and Quade [1968].

3.3 An analysis using "ridits"

This example is based on Schotz's [1966] study of lone drivers in injury producing accidents. The relevant data which have also been analyzed by Bhapkar and Koch [1968a] are shown in Table 3.

Table 3

Driver Group	Accident Type	Accident Severity				Total
		Minor	Moderate	Severe	Extreme	
Lone Driver	Rollover	21	567	1356	644	2588
Lone Driver	Non-Rollover	996	5454	2773	1256	10479
Injured Driver with passengers	Rollover	18	553	1734	869	3174
Injured Driver with passengers	Non-Rollover	679	4561	2516	1092	8848
Total		1714	11135	8379	3861	25089

One way of investigating the variable accident severity is to assign scores to the categories "minor," "moderate," "severe," and "extreme" and then to analyze the average values corresponding to the driver group \times accident type categories. A useful method for such scaling, which has been suggested by Bross [1958], is based on "ridits." Here these scores will be derived from the cumulative distribution of the ordinally valued variable accident severity in the combined total population. The methodology of section 2 applies to "ridits" because it enables both their direct calculation as well as their analysis in terms of driver group and accident type effects. Other aspects of their use are discussed by Bross [1958].

Let us now view the data as coming from a single multinomial population. Let \underline{p} denote the vector defined in (3.22). The \underline{A} -matrix in (3.23) produces all $\underline{p}' = [21, 567, 1356, 644, 996, 5454, 2773, 1256, 18, 553, 1734, 869, 679, 4561, 2516, 1092]/25089$ (3.22)

cell probabilities, all marginal probabilities for driver group \times accident type combinations, and the "ridit" scores. Each cell probability is then multiplied by the appropriate "ridit" and divided by the appropriate marginal probability on the log-scale with the \tilde{K} -matrix in (3.24). Finally, the average "ridit"

$$\tilde{A} = \begin{matrix} \tilde{I}_{16} \\ 24 \times 16 \\ \tilde{A}_1 \\ 8 \times 16 \end{matrix} \quad \text{where } \tilde{A}_1 = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 \\ 1 & \frac{1}{2} & 0 & 0 & 1 & \frac{1}{2} & 0 & 0 & 1 & \frac{1}{2} & 0 & 0 & 1 & \frac{1}{2} & 0 & 0 \\ 1 & 1 & \frac{1}{2} & 0 & 1 & 1 & \frac{1}{2} & 0 & 1 & 1 & \frac{1}{2} & 0 & 1 & 1 & \frac{1}{2} & 0 \\ 1 & 1 & 1 & \frac{1}{2} & 1 & 1 & 1 & \frac{1}{2} & 1 & 1 & 1 & \frac{1}{2} & 1 & 1 & 1 & \frac{1}{2} \end{bmatrix}$$

scores are derived by performing the additions corresponding to the \tilde{Q} -matrix in (3.25).

$$\tilde{K} = \begin{matrix} \tilde{I}_{16} \\ 16 \times 24 \\ \tilde{K}_1 \\ 16 \times 8 \end{matrix} \quad \text{where } \tilde{K}_1 = \begin{bmatrix} -\tilde{J}_4 & \tilde{O}_4 & \tilde{O}_4 & \tilde{O}_4 & \tilde{I}_4 \\ \tilde{O}_4 & -\tilde{J}_4 & \tilde{O}_4 & \tilde{O}_4 & \tilde{I}_4 \\ \tilde{O}_4 & \tilde{O}_4 & -\tilde{J}_4 & \tilde{O}_4 & \tilde{I}_4 \\ \tilde{O}_4 & \tilde{O}_4 & \tilde{O}_4 & -\tilde{J}_4 & \tilde{I}_4 \end{bmatrix} \quad (3.24)$$

$$\tilde{Q} = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{bmatrix} \quad (3.25)$$

The resulting \tilde{F} vector obtained from (2.10) with \tilde{A} , \tilde{K} , and \tilde{Q} as defined in (3.23), (3.24) and (3.25) is given in (3.26) together with the corresponding estimated covariance matrix $\tilde{V}_{\tilde{F}}$. The effects of driver group and accident type on the "ridit" score can be evaluated by fitting the linear model in (3.27). Where β_0 is an overall

$$\tilde{F} = \begin{bmatrix} 0.649 \\ 0.445 \\ 0.674 \\ 0.459 \end{bmatrix} \quad \text{and} \quad \tilde{V}_F = \begin{bmatrix} 18.10 & -1.82 & -1.21 & -1.79 \\ & 3.96 & -1.60 & -2.68 \\ & & 13.12 & -1.56 \\ & & & 5.18 \end{bmatrix} \times 10^{-6} \quad (3.26)$$

$$\varepsilon\{\tilde{F}\} = \tilde{X}\tilde{\beta} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & -1 \\ 1 & -1 & 1 \\ 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \quad (3.27)$$

mean, β_1 is a driver group effect, and β_2 is an accident type effect. The goodness of fit statistic in (2.4) corresponding to the model (3.27) is $X^2 = 2.31$ with D.F. = 1 which is non-significant ($\alpha = 0.10$). This result implies that there is no interaction between driver group and accident type with respect to the ridit score. The estimate \tilde{b} of $\tilde{\beta}$ in (3.27) and its estimated covariance matrix $\tilde{V}_{\tilde{b}}$ are shown in (3.28). Hence, it can be verified that the test

$$\tilde{b} = \begin{bmatrix} 0.557 \\ -0.009 \\ 0.105 \end{bmatrix} \quad \text{and} \quad \tilde{V}_{\tilde{b}} = \begin{bmatrix} 1.15 & 0.06 & 1.53 \\ & 2.53 & 0.30 \\ & & 2.86 \end{bmatrix} \times 10^{-6} \quad (3.28)$$

statistics for the hypotheses $H_0: \beta_1 = 0$ and $H_0: \beta_2 = 0$ are $X^2 = 31.18$ with D.F.=1 and $X^2 = 3874.29$ with D.F.=1 respectively, both of which are significant ($\alpha = .01$).

The analysis given here is different from the more common use of "ridits" as known constant coefficients rather than scores derived from the observed data. For that approach the data for the four rows of Table 3 are usually viewed as coming from separate and independent multinomial populations. In this event, the value of \tilde{F} is the same as in (3.26); but \tilde{V}_F is now the diagonal matrix shown in (3.29). Thus, the goodness of fit test for (3.27) is $X^2 = 2.29$ with D.F. = 1. Also, the estimate \tilde{b} is essentially the same as that shown in

$$\tilde{V}_F = \begin{bmatrix} 19.49 & & & \\ & 6.72 & & \\ & & 14.17 & \\ & & & 7.79 \end{bmatrix} \times 10^{-6} \quad (3.29)$$

(3.28), but the test statistics for $H_0 : \beta_1 = 0$ and $H_0 : \beta_2 = 0$ are slightly smaller with values of $X^2 = 31.17$ and $X^2 = 3707.00$ respectively. For the case in which there are several independent multinomial populations, this approach has some appeal; but since the "ridit" scores used here have still been determined from the observed data, some questions arise as to whether the components of \tilde{F} can be viewed as independent. This difficulty can, however, be avoided by using the methods of Section 2 in a somewhat different manner than what was originally described for these data. In particular, the appropriate \tilde{A} is given in (3.30) with $n_1 = 2588$, $n_2 = 10479$, $n_3 = 3174$, $n_4 = 8848$, and $n = 25089$;

$$\tilde{A} = \begin{bmatrix} I_{\sim 16} \\ A_{\sim 1} \\ 4 \times 16 \end{bmatrix} \quad \text{where } \tilde{A}_{\sim 1} = \frac{1}{n} \begin{bmatrix} \frac{1}{2}n_1 & 0 & 0 & 0 & \frac{1}{2}n_2 & 0 & 0 & 0 & \frac{1}{2}n_3 & 0 & 0 & 0 & \frac{1}{2}n_4 & 0 & 0 & 0 \\ n_1 & \frac{1}{2}n_1 & 0 & 0 & n_2 & \frac{1}{2}n_2 & 0 & 0 & n_3 & \frac{1}{2}n_3 & 0 & 0 & n_4 & \frac{1}{2}n_4 & 0 & 0 \\ n_1 & n_1 & \frac{1}{2}n_1 & 0 & n_2 & n_2 & \frac{1}{2}n_2 & 0 & n_3 & n_3 & \frac{1}{2}n_3 & 0 & n_4 & n_4 & \frac{1}{2}n_4 & 0 \\ n_1 & n_1 & n_1 & \frac{1}{2}n_1 & n_2 & n_2 & n_2 & \frac{1}{2}n_2 & n_3 & n_3 & n_3 & \frac{1}{2}n_3 & n_4 & n_4 & n_4 & \frac{1}{2}n_4 \end{bmatrix} \quad (3.30)$$

Similarly, the appropriate \tilde{K} is given in (3.31). If the \tilde{Q} in (3.25) were used,

$$\tilde{K} = \begin{bmatrix} I_{\sim 16} & K_{\sim 1} \end{bmatrix} \quad \text{where } K_{\sim 1} = \begin{bmatrix} I_{\sim 4} \\ I_{\sim 4} \\ I_{\sim 4} \\ I_{\sim 4} \end{bmatrix} \quad (3.31)$$

then the resulting $\tilde{V}_{\tilde{F}}$ would be singular since determination of the ridits in this situation removes one of the degrees of freedom associated with the four driver group x accident type combinations. Hence, to obtain results analogous to those previously given, one can use the \tilde{Q} matrix in (3.32).

$$\tilde{Q} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 \\ 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 \end{bmatrix} \quad (3.32)$$

The resulting function vector \tilde{F} obtained from (2.10) with \tilde{A} , \tilde{K} , and \tilde{Q} as defined in (3.30), (3.31), and (3.32) is given in (3.33) together with the corresponding estimated covariance matrix $\tilde{V}_{\tilde{F}}$

$$\tilde{F} = \begin{bmatrix} -0.040 \\ 0.420 \\ -0.011 \end{bmatrix} \quad \text{and} \quad \tilde{V}_{\tilde{F}} = \begin{bmatrix} 48.15 & 6.57 & 19.15 \\ & 46.11 & 4.28 \\ & & 48.15 \end{bmatrix} \times 10^{-6} \quad (3.33)$$

From these results, it follows that test statistics for driver group effects, accident type effects, and their interaction are $X^2 = 32.73$, $X^2 = 3825.63$, and $X^2 = 2.29$ respectively, each with D.F. = 1.

Finally, other types of ridit analysis like that of Williams and Grizzle [1972] based on partial scores can be undertaken within this same framework.

3.4. An analysis of partial association

The data in Table 4 are from the Kihlberg, Narragon, and Campbell [1964] study of the relationship between car size and accident injuries for drivers who were alone at the time of an automobile accident. Here, we shall assume

Table 4

Car Weight	Accident Type	Driver Not Ejected		Driver Ejected	
		Not Severe	Severe	Not Severe	Severe
Small	Collision	350	150	26	23
Small	Roll-over	60	112	19	80
Standard	Collision	1878	1022	111	161
Standard	Roll-over	148	404	22	265

that each of the rows of Table 4 have arisen from separate and independent multinomial populations. In this case, one question of interest is the manner in which the partial association between accident severity and driver ejection after adjustment for car weight depends on accident type. Given that the relationship between accident severity and ejection is similar for the two car weights, this

problem can be dealt with by considering for each accident type, an index which is the ratio of the observed number of severe accidents where the driver was ejected to the expected number of such accidents when accident severity and driver ejection are independent. This index, which has been also extensively discussed by Campbell [1970], can be analyzed by the methods described in Section 2.

The \tilde{A} -matrix in (3.30) produces the cell probability of a severe accident with driver ejected together with the corresponding marginals for each car weight \times accident type combination. The \tilde{K} -matrix in (3.31) then produces the appropriate

$$\tilde{A} = \begin{matrix} 12 \times 16 \\ \left[\begin{array}{cccc} \tilde{A}_1 & & & \\ & \tilde{A}_1 & & \\ & & \tilde{A}_1 & \\ & & & \tilde{A}_1 \end{array} \right] \end{matrix} \quad \text{where} \quad \tilde{A}_1 = \begin{matrix} 3 \times 4 \\ \left[\begin{array}{cccc} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{array} \right] \end{matrix} \quad (3.30)$$

observed and expected proportions on the log-scale. These observed and expected

$$\tilde{K} = \begin{matrix} 8 \times 12 \\ \left[\begin{array}{cccc} \tilde{K}_1 & & & \\ & \tilde{K}_1 & & \\ & & \tilde{K}_1 & \\ & & & \tilde{K}_1 \end{array} \right] \end{matrix} \quad \text{where} \quad \tilde{K}_1 = \begin{matrix} \left[\begin{array}{cc} 1 & 0 \\ 0 & 1 \end{array} \right] \end{matrix} \quad (3.31)$$

values are then added across car weight within each accident type by use of the \tilde{Q} -matrix in (3.32). The index representing the ratio of observed over expected

$$\tilde{Q} = \begin{matrix} 4 \times 8 \\ \left[\begin{array}{ccccccccc} 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{array} \right] \end{matrix} \quad (3.32)$$

for each accident type is obtained on the log-scale with the \tilde{L} matrix in (3.33). For the data in Table 4, the function vector \tilde{F} obtained from expression (2.12)

$$\tilde{L} = \begin{matrix} \left[\begin{array}{cc} 1 & -1 \\ 0 & 0 \end{array} \right] \quad \left[\begin{array}{cc} 0 & 0 \\ 1 & -1 \end{array} \right] \end{matrix} \quad (3.33)$$

with \tilde{A} , \tilde{K} , \tilde{Q} , and \tilde{L} defined by (3.30), (3.31), (3.32), and (3.33) is given in (3.34) together with its estimated covariance matrix \tilde{V}_F .

$$\tilde{F} = \begin{bmatrix} 0.433 \\ 0.139 \end{bmatrix} \quad \tilde{V}_F = \begin{bmatrix} 50.11 \\ 5.47 \end{bmatrix} \times 10^{-4} \quad (3.34)$$

If the linear model in (3.35) is fitted to \tilde{F} in (3.34) then the goodness of fit

$$\varepsilon\{\tilde{F}\} = \tilde{X}\tilde{\beta} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \beta_0 \quad (3.35)$$

test indicates whether the partial association index between accident severity and driver ejection adjusted for car weight is the same for the two accident types. The appropriate test statistic obtained from (2.4) is $X^2 = 15.49$ with D.F. = 1 which is significant ($\alpha = .01$). Hence, the relationship of severity to driver ejection is different for collision and roll-over accidents.

In addition, it is of interest to test the significance of each of the indices in \tilde{F} . This can be done by using the linear model in (3.36) and testing the hypotheses $H_0: \beta_1 = 0$ and $H_0: \beta_2 = 0$. The results based on (2.5) are $X^2 = 37.42$ and $X^2 = 35.32$ respectively, each of which is significant ($\alpha = .01$).

$$\varepsilon\{\tilde{F}\} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \quad (3.36)$$

For the two accident types, each of these test statistics is analogous to what would be obtained using the method of Mantel and Haenszel [1959] or Mantel [1963] except that they are based on the log-scale functions F with estimated covariance matrix obtained from (2.13). The advantage of the approach suggested here is that it permits the investigation of other models for indices of partial association like that in (3.35).

An alternative approach to this problem, is to use \tilde{K} in (3.37) to calculate the observed over expected ratio on the log scale for each car weight \times accident

$$\tilde{K} = \begin{bmatrix} 1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 & -1 \end{bmatrix} \quad (3.37)$$

type combination. The resulting function vector \tilde{F} (there is no \tilde{Q} or \tilde{L} matrix here), is then considered in terms of the linear model in (3.38) where β_0 is an overall mean, β_1 is a car weight effect, and β_2 is an accident type effect.

$$\tilde{\epsilon}\{\tilde{F}\} = \tilde{X}\tilde{\beta} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & -1 \\ 1 & -1 & 1 \\ 1 & -1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} \quad (3.38)$$

The goodness of fit statistic from (2.4) is $X^2 = 0.09$ with D.F. = 1. Also, test statistics for $H_0: \beta_1 = 0$ and $H_0: \beta_2 = 0$ from (2.5) are $X^2 = 0.18$ with D.F. = 1 and $X^2 = 39.33$ with D.F. = 1 respectively. Since the goodness of fit statistic and the statistic for $H_0: \beta_1 = 0$ are non-significant ($\alpha = .10$), it follows that this measure of association between accident severity and driver ejection does not depend on car weight. Thus, this assumption which was required by our initial analysis in (3.30) - (3.36) is appropriate for these data. As a result, the test statistic for $H_0: \beta_2 = 0$ and the goodness of fit statistic for (3.35) are both directed at the effect of accident type; and it is useful to note that they have similar values. The advantage of this analysis is that it allows the association between accident severity and ejection to be investigated at the more refined level corresponding to all car weight \times accident type combinations. However, in many situations of this kind, the sample sizes are sometimes too small for such an extensive analysis. Hence, the approach which was emphasized at the beginning of this section may be more feasible and appropriate.

Finally, one can note that the results obtained here are somewhat different

from those obtained by Bhapkar and Koch [1968a] in a previous consideration of these data. They used the log - cross-product-ratio measure of association between accident severity and driver ejection and, on applying a model similar to (3.38), found that it was not significantly ($\alpha = .05$) affected by car weight and accident type. Thus, the log-observed-over-expected-ratio for the severe accident with driver ejected category is a more sensitive index to accident type effects than the log-cross-product-ratio in this application. This statement, however, should be interpreted carefully because use of the log-observed-over-expected ratio for a given category is justified only when the corresponding cell in the 2×2 table can be viewed as having the greatest practical importance.

REFERENCES

- Berkson, J. [1955]: Maximum likelihood and minimum χ^2 estimates of the logistic function. J. Amer. Statist. Ass. 50, 130-62.
- Bhapkar, V.P. and Koch, G.G. [1968a]. Hypotheses of 'no interaction' in multidimensional contingency tables. Technometrics 10, 107-23.
- Bhapkar, V.P. and Koch, G.G. [1968b]. On the hypotheses of 'no interaction' in contingency tables. Biometrics 24, 567-94.
- Bishop, Y.M.M. [1969]. Full contingency tables, logits, and split contingency tables. Biometrics 25, 383-99.
- Bishop, Y.M.M. [1971]. Effects of collapsing multidimensional contingency tables. Biometrics 27, 545-62.
- Bross, I.D.J. [1958]. How to use 'ridit' analysis. Biometrics 14, 18-38.
- Campbell, B.J. [1970]. Driver injury in automobile accidents involving certain car models. University of North Carolina Highway Safety Research Center Report.
- Davis, C.E. and Quade, D. [1968]. "On comparing the correlations within two pairs of variables," Biometrics 24, 987-996.
- Goodman, L.A. [1970]. The multivariate analysis of qualitative data: interactions among multiple classifications. J. Amer. Statist. Ass. 65, 226-56.
- Goodman, L.A. [1971a]. The partitioning of chi-square, the analysis of marginal contingency tables, and the estimation of expected frequencies in multidimensional contingency tables. J. Amer. Statist. Ass. 66, 339-44.
- Goodman, L.A. [1971b]. The analysis of multidimensional contingency tables; Stepwise procedures and direct estimation methods for building models for multiple classifications. Technometrics 13, 33-61.
- Goodman, L.A. and Kruskal, W.H. [1954]. Measures of association for cross classifications. J. Amer. Statist. Ass. 49, 734-64.
- Goodman, L.A. and Kruskal, W.H. [1959]. Measures of association for cross classifications. II: Further discussion and references. J. Amer. Statist. Ass. 54, 123-63.
- Goodman, L.A. and Kruskal, W.H. [1963]. Measures of association for cross classification. III: Approximate sampling theory. J. Amer. Statist. Ass. 58, 310-64.
- Grizzle, J.E., Starmer, C.F., and Koch, C.G. [1969]. Analysis of categorical data by linear models. Biometrics 25, 489-504.
- Kihlberg, J.K., Narragon, E.A., and Campbell, B.J. [1964]. Automobile crash injury in relation to car size. Cornell Aeronautical Laboratory Report No. V.J.-1823-R11.
- Ku, H.H., Varner, R., and Kullback, S. [1971]. Analysis of multidimensional contingency tables, J. Amer. Statist. Ass. 66, 55-64.

- Mantel, N. and Haenszel, W. [1959]. Statistical aspects of the analysis of data from retrospective studies of disease. J. Nat. Cancer Inst. 22, 719-48.
- Mantel, N. [1963]. Chi-square tests with one degree of freedom; extensions of the Mantel-Haenszel procedure. J. Amer. Statist. Ass. 58, 690-700.
- Neyman, J. [1949]. Contribution to the theory of the χ^2 test. Pp. 239-73 in Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability, University of California Press, Berkeley and Los Angeles.
- Roy, S.N. and Kastenbaum, M.A. [1956]. On the hypothesis of no interaction in a multiway contingency table. Ann. Math. Statist. 27, 749-57.
- Schotz, W.E. [1966]. The lone driver involved in injury-producing accidents. Cornell Aeronautical Laboratory Report No. VJ-1823-R19.
- Wald, A. [1943]. Tests of statistical hypotheses concerning several parameters when the number of observations is large. Trans. Amer. Math. Soc. 54, 426-82.
- Williams, O. D. and Grizzle, J. E. [1972]. Analysis of contingency tables having ordered response categories. J. Amer. Statist. Ass. 66.