

APPLICATIONS OF WEIGHTED LEAST SQUARES METHODS FOR  
FITTING VARIATIONAL MODELS TO HEALTH SURVEY DATA

Gary G. Koch and Maura E. Stokes, University of North Carolina, Chapel Hill

and

Dwight Brock, National Center for Health Statistics

1. Introduction

Data from complex sample surveys can be analyzed by using weighted least squares (WLS) methods similar to those described by Grizzle, Starmer and Koch [1969] for the analysis of categorical data. This approach allows the variation among domain estimates to be investigated using linear regression model strategies, provided that such estimates can be presumed to have an approximate multivariate normal distribution as a consequence of large sample size considerations. Such models can be formulated as orthocomplement matrices to constraint matrices for hypotheses with which the variation among the domain estimates are compatible, such hypotheses corresponding to sources of variation which are essentially equivalent to sampling variability. Specifically, this paper summarizes two examples of this type of analysis for domain estimates from the First Health and Nutrition Examination Survey (HANES) which was conducted during 1971-1974. One of these is concerned with percentage estimates of extreme armgirth (>40 cm) for an age x sex cross-classification; and the other is concerned with percentage estimates of persons having a regular dentist for an age x sex x income cross-classification. Both of these sets of estimates were obtained by combining the observed data for the subjects in this national probability sample in an appropriate way with respect to the survey design to produce percentage estimates for the United States target population. Post-stratification was used in order to adjust for the oversampling components of the HANES design for pre-school children, women of child-bearing age, elderly people and low income people, such oversampling having been undertaken so that the survey would provide more reliable estimates for these subpopulations.

For these two examples, the basic steps in the analysis of variation among the domain estimates are described. Also, attention is given to statistical issues concerning the evaluation of model goodness of fit and the use of model predicted values of domain estimates for inferential purposes.

2. Methodology

The vector  $\tilde{F}$  of extreme armgirth estimates and its corresponding covariance matrix  $\tilde{V}_F$  are shown in Table 1. The covariances here were calculated according to the method of balanced repeated replication described in McCarthy [1969] and Kish and Frankel [1970]. When the vector  $\tilde{F}$  of estimates is constructed from large samples like those in HANES, the estimates have an approximate normal distribution and linear hypotheses of the form

$$H_0: \tilde{W}\tilde{F} \hat{=} \tilde{Q},$$

concerning age, sex, and age x sex interaction can be tested to assess variation. (Here,  $\tilde{W}$  is a full rank matrix of contrast constraints.) These hypotheses are tested using the Wald statistic (quadratic form)

$$Q_{W,C} = \tilde{F}'\tilde{W}'(\tilde{W}\tilde{V}_F\tilde{W}')^{-1}\tilde{W}\tilde{F},$$

which has an approximate chi-square distribution with D.F.=Rank( $\tilde{W}$ ) under the null hypothesis;  $\tilde{V}_F$  represents the large sample (consistent) estimate of the covariance structure for  $\tilde{F}$ .

The hypothesis  $H_0$  can be interpreted as a goodness of fit test for the variational model  $\tilde{F} \hat{=} \tilde{X}\tilde{b}$ , implied when the hypothesis is accepted where  $\tilde{X}$  is a design matrix orthogonal to  $\tilde{W}$  and  $\tilde{b}$  is a vector of estimated parameters. In other words,  $\tilde{W}\tilde{F} \hat{=} \tilde{W}\tilde{X}\tilde{b} = \tilde{0}$  implies  $\tilde{F} \hat{=} \tilde{X}\tilde{b}$ .  $\tilde{W}$  is called the constraint formulation of the model and  $\tilde{X}$  is called the model specification formulation or the freedom equation formulation, which characterizes  $\tilde{F}$  in a manner compatible with  $H_0$ . If the goodness of fit for the model  $\tilde{X}$  is considered adequate, then WLS methods can be used to obtain the estimates  $\tilde{b} = (\tilde{X}'\tilde{V}_F^{-1}\tilde{X})^{-1}\tilde{X}'\tilde{V}_F^{-1}\tilde{F}$  and its estimated covariance matrix  $\tilde{V}_{\tilde{b}} = (\tilde{X}'\tilde{V}_F^{-1}\tilde{X})^{-1}$ . Since  $\tilde{b}$  will also be multivariate normal for large sample sizes, the Wald statistic

$$Q_{W,C} = \tilde{b}'\tilde{C}'(\tilde{C}\tilde{V}_{\tilde{b}}\tilde{C}')^{-1}\tilde{C}\tilde{b}$$

can be used to test hypotheses of the form  $H_{C,W}: \tilde{C}\tilde{b} \hat{=} \tilde{Q}$ .  $Q_{W,C}$  has an approximate chi-square distribution under the null hypothesis with D.F.=Rank( $\tilde{C}$ ).

3. Results for Extreme Armgirth Data Analysis

In Table 2, the preliminary hypotheses which were investigated are shown together with the corresponding constraint matrices and resulting chi-square test statistics.

The test statistics for  $H_1-H_3$  are significant ( $\alpha=.01$ ) and thus contradict the respective hypotheses. The additional hypothesis ( $H_4$ ) considered was that of no age x sex interaction. Since the corresponding Wald statistic of 4.73 (D.F.=4) is non-significant with  $p > .25$ , the hypothesis  $H_4$  is considered to be compatible with the estimates at hand. Thus, the variation among the estimates can be characterized by a linear model  $\tilde{F} = \tilde{X}_1\tilde{b}$ , where  $\tilde{X}_1$  is an orthocomplement to  $\tilde{W}_4$ . In addition, the goodness of fit chi-square for this model will be identical to the Wald test statistic for the hypothesis  $H_4$ . The specification matrix for the model  $\tilde{X}_1$  is given below together with estimates for its parameters.

$\tilde{F} = \tilde{X}_1\tilde{b} =$	$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}$	$\begin{cases} \text{reference values for} \\ \text{male 25-34} \\ \text{increment for female} \\ \text{increment for age 35-44} \\ \text{increment for age 45-54} \\ \text{increment for age 55-64} \\ \text{increment for age 65-74} \end{cases}$
$\tilde{b} =$	$\begin{bmatrix} 1.30 \\ 1.58 \\ 1.52 \\ -.23 \\ -.37 \\ -1.13 \end{bmatrix}$	$\begin{cases} \text{Goodness of fit} \\ \text{statistic (D.F.=4)} = 4.73 \\ (p = 0.32) \end{cases}$

Since the variation among the domain estimates is adequately characterized by the model  $X_1$ , further analyses can be based on the resulting estimate  $\hat{b}$ . More specifically, hypotheses of the form  $H_0: C\hat{b} = 0$  can be tested by chi-square statistics of the type  $Q_{W,C} = \hat{b}'C'(CV_{\hat{b}} C')^{-1}C\hat{b}$  where  $V_{\hat{b}}$  is the covariance matrix for  $\hat{b}$ . In Table 3, some hypotheses concerning  $\hat{b}$  are shown together with the corresponding contrast matrices and the resulting test statistic. Clearly, the hypotheses  $H_5$  and  $H_6$  are contradicted by the data ( $\alpha=.01$ ); and  $H_7$  is judged to be compatible with the data. This implies that the variation among the domain estimates can be characterized by a lower dimensional model  $X_2$ , where  $X_2$  is as follows:

$$X_2 B_2 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \text{reference value for} \\ \text{males, 25-34} \\ \text{increment for} \\ \text{females} \\ \text{increment for age} \\ \text{65-74} \end{bmatrix}$$

The goodness of fit test statistic for  $X_2$  is 6.07 (with D.F.=7), for which the p-value is 0.53. This result can be interpreted as the goodness of fit statistic for  $X_1$  incremented by an amount equal to that of the Wald statistic for  $H_7$  since  $H_7$  implies  $\hat{b} = Zg$ , which in combination with  $H_4$ , implies  $F = X_1 Zg$  which is equivalent to  $W_{C,F} = 0$  where  $W_C$  is orthogonal to  $X_1 Z = X_2$ . So  $Q_7$  and  $Q_4(X_1)$  represent additive components for the goodness of fit test statistic of  $X_2$ . Table 4 contains the parameter estimates for the model  $X_2$ . These can be interpreted as indicating that the percentage extreme armgirth estimates were 1.68 higher for females than males for all age domains, and .95 lower for ages 65-74 in both sex domains than the other age ranges. Also, Table 4 includes the model predicted estimates as well as their standard errors. Also, it can be noticed there that the predicted estimates have smaller estimated standard errors than the original estimates because they are based on the combined information for all domains through its three estimated parameters instead of the separate information for the individual domains.

It should be pointed out that the additive model  $F = X_1 \hat{b}$  was an entirely adequate model for the cross-classified estimates of extreme armgirth percentages, as indicated by its goodness of fit test. The use of the model  $F = X_2 \hat{b}$  and the resulting decrease in the number of parameters to be estimated allows a more simplified model for the domain estimates, but it may not necessarily be considered the best one. In fact, a test of the linearity of age effect of the form  $C\hat{F} = 0$  where

$$C = \begin{bmatrix} 0 & 0 & -2 & 1 & 0 & 0 \\ 0 & 0 & -3 & 0 & 1 & 0 \\ 0 & 0 & -4 & 0 & 0 & 1 \end{bmatrix}$$

has a Wald statistic of 3.93 with D.F.=3 and

$p = 0.27$ . It corresponds to a variational model

$$X_3 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 & 4 & 0 & 1 & 2 & 3 & 4 \end{bmatrix}$$

and thus presents a paradox by implying variation among the age  $\leq 64$  subdomains. Thus, choice of a model is not solely dictated by the desire for one which is reduced to as low a dimension as possible, but rather, by the desire to have one which is both as parsimonious as possible and whose interpretation is reasonable for the data at hand.

Thus, the model  $X_3$  may have been appropriate if there were an a priori basis for its consideration. In such a situation, there would have been no need for investigative hypothesis testing of the type that included hypothesis  $H_7$ ; the design matrix  $X_3$

would have been fitted immediately and its goodness of fit examined to assess if the fit was indeed satisfactory. At this point, it should be recognized that the model  $X_2$  was deduced in an a posteriori fashion.

For this reason, if the objective of analysis was to detect significant sources of variation and make inferences involving the resulting model estimates, then multiple comparison type approaches would need to be taken into account in order to assess significance and to derive analogous confidence intervals. In this regard, either Bonferroni inequality, Scheffe type methods or their combination can be used with the relevant consideration being the formulation of the range of hypotheses which are of interest and the types of models which could be regarded as having an a priori basis even if they were not so specified. For example, the no interaction model  $X_2$  might be

considered plausible on a priori grounds and so simultaneous inference with respect to its parameters could be undertaken by Scheffe methods relative to chi-square approximations with D.F.=Rank( $X_2$ )=6. Also, it can be noted that such multiple comparison methods are applicable to confidence intervals for predicted values as well as to tests of significance and thereby represent a strategy for dealing with dilemmas concerning model overfitting in a posteriori situations.

Finally, sometimes hypothesis testing and model fitting are undertaken purely for exploratory purposes with respect to assessing the relative extent of different sources of variation. In these cases, inferences in a technical sense are not an objective of analysis because a multiplicity of descriptive interpretations may be indicated as plausible. Thus, the use of multiple comparison procedures or other analysis strategies with a similarly oriented inferential spirit may not be necessary, provided that interpretations of results are suitably qualified. For further discussion of this example, see Koch and Stokes [1979]; and of related statistical issues, see Koch, Gillings, and Stokes [1980] and Koch and Stokes [1981].

4. Results for Regular Dentist Data Analysis

The analysis of the estimates involving the attribute of having a regular dentist is somewhat more complicated than that of the armgirth estimates because it involves a three-way cross-classification for age, sex, and income. The specific cross-classification scheme under investigation,

the estimates for regular dentist percentages, and their standard errors are given in Table 5.

For the armgirth data, the process for determining an adequate model involved assessing most of the pertinent sources of variation in domain estimates and then fitting the model which the results implied. A three-way cross-classification involves consideration of many more potential sources of variation, and so the model fitting process requires a relatively practical strategy for screening potential models. Thus, the first step for analysis in such situations is to fit several relatively simple models involving individual sources of variation and partially additive combinations of them in order to assess their goodness of fit rather than to test hypotheses of the constraint form  $H_0: \underline{W}F \hat{=} 0$  as outlined in Section 3. However, the results of such tests are interpreted in the same spirit of those for their constraint formulation counterparts with respect to the identification of a preliminary model as a framework for further analysis. The model specification matrices which were used for this preliminary purpose are shown in Table 6 together with the corresponding hypothesis descriptions, goodness of fit test statistics and p-values.

Clearly, the hypotheses  $H_1, H_2, H_3$  are contradicted, and so the variability among the age x sex subdomains of the age domains, and the age x income subdomains of the sex domains are greater than that which would be expected with respect to their inherent sampling variability. Also, this same conclusion holds if the significance of these test statistics is evaluated from a Scheffe multiple comparison point of view by reference to a chi-square distribution with 39 degrees of freedom (i.e., the dimension of the overall vector space of comparison contrasts among domain estimates of which  $H_1, H_2, H_3$  are subsets).

An hypothesis which is indicated to be compatible with the variation among the domain estimates is the one corresponding to no interaction among the age, sex, and income sources in the sense of the additive model  $X_4$ . The goodness of fit test statistic for this model is  $Q(X_4) = 34.85$  (D.F.=31 and  $p=0.29$ ). The estimated values for the parameters of this model and their estimated standard errors are shown below:

Reference value for males 25-34, < \$3000	35.99 + 3.04
Increment for females	11.20 + 1.26
Increment for age 35-44	6.74 + 1.45
Increment for age 45-54	1.69 + 1.51
Increment for age 55-64	-3.77 + 2.01
Increment for age 65-74	-8.21 + 2.45
Increment for income \$3-9999	19.08 + 2.51
Increment for income \$10-19999	34.72 + 2.30
Increment for income > \$20000	43.92 + 2.36

As stated previously, statistical tests to demonstrate no interaction among the age, sex, and income sources of variation could have been undertaken in more detail via hypothesis testing of the type  $H_0: \underline{W}F \hat{=} 0$ . In this case the constraints  $\underline{W}$  would have been constructed to represent appropriate difference functions which would be zero under null hypotheses for various formulations of no second or third order interaction. While good

practice, this process can be tedious to implement computationally since it involves rather cumbersome matrices. For this reason, details concerning such tests are not included here; they are documented in Koch and Stokes [1979].

In Table 7, results for tests of hypotheses concerning the parameter vector estimated for the model  $X_4$  are given. These hypotheses have the general form  $H_0: \underline{C}\underline{b} \hat{=} 0$  where  $\underline{C}$  is the corresponding hypothesis specification matrix. All of these hypotheses are clearly contradicted at the  $\alpha=0.01$  significance level; they are also contradicted if each is evaluated from the Scheffe multiple comparison point of view by reference to a chi-square approximation with D.F.=39 (as discussed previously). These results suggest that further simplification of the model  $X_4$  may not be warranted (without some type of a priori justification) since all of the sources of variation corresponding to its parameters are significant. Thus,  $\underline{F} \hat{=} \underline{X}_4 \underline{b}$  is considered an adequate characterization of the variation among the regular dentist percentage estimates for the respective domains.

Finally, the predicted values obtained by using the model  $X_4$  and their standard errors are given in

Table 8. These indicate that the percentages of persons with a regular dentist is larger for females than males, increases with income, and initially increases with age to the 35-44 year range and then decreases. Thus, for example, the lowest predicted value (27.8%) corresponds to males in the 65-74 years age range and the less than \$3000 income range; and the highest (97.9%) corresponds to females in the 35-44 years age range and the  $\geq$  \$20,000 income range.

#### REFERENCES

- Grizzle, J.E., Starmer, C.F., and Koch, G.G. (1969). Analysis of categorical data by linear models. *Biometrics* 25, 489-504.
- Kish, L. and Frankel, M. (1970). Balanced repeated replications for standard errors. *Journal of the American Statistical Association* 65, 1071-1094.
- Koch, G.G., and Bhapkar, V.P. (1981). Chi-square tests. Entry to appear in *Encyclopedia of Statistical Sciences*, Norman L. Johnson and Samuel Kotz, eds., to be published by J. Wiley & Sons, Inc.
- 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.
- Koch, G.G., Gillings, D.B., and Stokes, M.E. (1980). Biostatistical implications of design, sampling and measurement to the analysis of health science data. *Annual Review of Public Health* 1, 163-225.
- Koch, G.G. and Stokes, M.E. (1979). Annotated computer applications of weighted least squares methods for illustrative analyses of examples involving health survey data. Technical report prepared for U.S. National Center for Health Statistics.
- Koch, G.G. and Stokes, M.E. (1981). Chi-square tests: numerical examples. Entry to appear in *Encyclopedia of Statistical Sciences*, Norman L. Johnson and Samuel Kotz, eds., to be published by John Wiley & Sons, Inc.
- McCarthy, P.J. (1969). Pseudoreplication: Half-samples. *International Statistical Review*, 37, 239-264.
- United States National Center for Health Statistics (1973). Plan and operation of the health and nutrition examination survey. *Vital and Health Statistics, Series 1, Numbers 10A, 10B, 14* DHEW Pub. No. (HRA), 73-1310.

TABLE 1: ARMGIRTH PERCENTAGE ESTIMATES AND ESTIMATED COVARIANCE MATRIX FROM THE 1971-1974 HANES SURVEY OF THE UNITED STATES POPULATION

DOMAINS		EXTREME ARMGIRTH PERCENTAGE ESTIMATES	BALANCED REPEATED REPLICATION ESTIMATED COVARIANCE MATRIX X 10 <sup>4</sup> FOR EXTREME ARMGIRTH PERCENTAGE ESTIMATES									
SEX	AGE		4529	1375	-45	-68	222	-591	235	-212	-693	-467
Male	25-34	2.14										
Male	35-44	2.08		5191	-423	392	291	-72	-39	269	-671	-448
Male	45-54	0.77			901	-203	-52	-234	192	-758	276	18
Male	55-64	0.79				1850	-185	-560	-594	83	-848	-223
Male	65-74	0.29					218	-26	153	-22	93	-57
Female	25-34	2.78						1703	202	181	855	-72
Female	35-44	2.88	SYMMETRIC						2515	-759	789	79
Female	45-54	3.46								3533	-770	-188
Female	55-64	2.44									3016	231
Female	65-74	1.55										1004

TABLE 2: LINEAR HYPOTHESES REGARDING ARMGIRTH ESTIMATES AND RESULTING TEST STATISTICS

HYPOTHESIS	CONTRAST MATRIX W	CHI-SQUARE STATISTIC	D.F.	P-VALUE
H <sub>1</sub> : There is no difference between the sex subdomains of each age domain.	$W_1: \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & -1 \end{bmatrix}$	35.26	5	< 0.001
H <sub>2</sub> : There is no variation among the age domains for males.	$W_2: \begin{bmatrix} 1 & 0 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$	14.97	4	0.005
H <sub>3</sub> : There is no variation among the age domains for females.	$W_3: \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{bmatrix}$	12.17	4	0.016
H <sub>4</sub> : There is no variation among the age domains for the differences between males and females	$W_4: \begin{bmatrix} 1 & -1 & 0 & 0 & 0 & -1 & 1 & 0 & 0 & 0 \\ 1 & 0 & -1 & 0 & 0 & -1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & -1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & -1 & -1 & 0 & 0 & 0 & 1 \end{bmatrix}$	4.75	4	0.316

TABLE 3: LINEAR HYPOTHESES CONCERNING THE PARAMETER ESTIMATES FOR THE MODEL X<sub>1</sub>

HYPOTHESIS	CONTRAST MATRIX	CHI-SQUARE STATISTIC	D.F.	P-VALUE
H <sub>5</sub> : There is no variation between male and female subdomains of age domains given X <sub>1</sub>	[0 1 0 0 0]	30.53	1	< 0.001
H <sub>6</sub> : There is no variation among age subdomains of the sex domains given X <sub>1</sub>	$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	23.89	4	< 0.001
H <sub>7</sub> : There is no variation among the age < 64 subdomains of the sex domains given X <sub>1</sub> .	$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$	1.34	3	0.719

TABLE 4: EXTREME ARMGIRTH PERCENTAGE ESTIMATES FOR THE 1971-1974 HANES SURVEY

Sex	Age	HANES Extreme Armgirth Percentage Estimate s.e.		Simplified Linear Model Structure (X)		Parameter Estimate + s.e.'s	Model Predicted Extreme Armgirth Estimates s.e.	
		Male	25-34	2.14	0.67	1	0	0
Male	35-44	2.08	0.72	1	0	0	1.08	0.17
Male	45-54	0.77	0.30	1	0	0	1.08	0.17
Male	55-64	0.79	0.43	1	0	0	1.08	0.17
Male	65-74	0.29	0.15	1	0	1	0.13	0.12
Female	25-34	2.78	0.41	1	1	0	2.76	0.18
Female	35-44	2.88	0.50	1	1	0	2.76	0.18
Female	45-54	3.46	0.59	1	1	0	2.76	0.18
Female	55-64	2.44	0.55	1	1	0	2.76	0.18
Female	65-74	1.55	0.32	1	1	1	1.81	0.26

TABLE 5: ESTIMATED REGULAR DENTIST PERCENTAGES + STANDARD ERRORS FOR THE 1971-1974 HANES SURVEY

Sex	Age	Income Classifications			
		< 3000	3000-9999	10000-19999	≥20000
M	25-34	30.3 + 18.0	56.3 + 6.3	67.6 + 4.5	75.7 + 7.9
M	35-44	30.2 + 14.7	60.9 + 6.9	75.7 + 5.4	88.0 + 5.2
M	45-54	39.6 + 10.9	53.8 + 5.1	78.4 + 4.0	85.6 + 5.3
M	55-64	28.9 + 10.4	46.6 + 6.5	63.7 + 5.9	81.0 + 9.2
M	65-74	28.4 + 6.4	45.2 + 4.3	51.4 + 10.0	86.5 + 10.5
F	25-34	37.1 + 14.4	70.1 + 4.6	83.4 + 3.2	87.1 + 6.0
F	35-44	53.0 + 11.8	69.6 + 6.5	84.8 + 3.4	90.5 + 5.5
F	45-54	56.3 + 12.4	65.6 + 5.2	79.6 + 3.6	90.4 + 4.2
F	55-64	39.4 + 7.7	58.8 + 5.7	80.6 + 4.9	92.0 + 5.7
F	65-74	42.1 + 7.0	63.5 + 4.7	63.7 + 11.1	64.8 + 12.4

TABLE 7: MODEL X<sub>4</sub> PREDICTED REGULAR DENTIST PERCENTAGES + STANDARD ERRORS FOR THE 1971-1974 HANES SURVEY

Sex	Age	Income Classifications			
		< 3000	3000-9999	10000-19999	>20000
M	25-34	36.0 + 3.0	55.1 + 1.7	70.7 + 1.7	79.9 + 1.8
M	35-44	42.7 + 2.9	61.8 + 2.0	77.5 + 1.8	86.7 + 1.5
M	45-54	37.7 + 2.6	56.8 + 1.5	72.4 + 1.8	81.6 + 1.4
M	55-64	32.2 + 2.7	51.3 + 2.1	67.0 + 2.0	76.1 + 2.0
M	65-74	27.8 + 2.8	46.9 + 1.4	62.5 + 2.1	71.7 + 2.2
F	25-34	47.2 + 3.1	66.3 + 1.8	81.9 + 1.6	91.1 + 1.9
F	35-44	54.0 + 2.9	73.0 + 1.9	88.7 + 1.6	97.9 + 1.5
F	45-54	48.9 + 2.6	68.0 + 1.4	83.6 + 1.6	92.8 + 1.4
F	55-64	43.4 + 2.5	62.5 + 1.8	78.1 + 1.6	87.3 + 1.7
F	65-74	39.0 + 3.1	58.1 + 1.7	73.7 + 2.3	83.0 + 2.5

TABLE 6: TEST STATISTICS FOR GOODNESS OF FIT OF LINEAR MODELS FOR REGULAR DENTIST DATA

SEX	POPULATION		MODEL 1 $X_1$	MODEL 2 $X_2$	MODEL 3 $X_3$	MODEL 4 $X_4$
	AGE (YEARS)	INCOME (\$1,000's)				
M	25-34	< 3	1 0	1 0 0 0 0	1 0 0 0	1 0 0 0 0 0 0 0
M	25-34	3-9	1 0	1 0 0 0 0	0 1 0 0	1 0 0 0 0 0 1 0 0
M	25-34	10-19	1 0	1 0 0 0 0	0 0 1 0	1 0 0 0 0 0 0 1 0
M	25-34	>20	1 0	1 0 0 0 0	0 0 0 1	1 0 0 0 0 0 0 0 1
M	35-44	< 3	1 0	0 1 0 0 0	1 0 0 0	1 0 1 0 0 0 0 0 0
M	35-44	3-9	1 0	0 1 0 0 0	0 1 0 0	1 0 1 0 0 0 0 1 0 0
M	35-44	10-19	1 0	0 1 0 0 0	0 0 1 0	1 0 1 0 0 0 0 0 1 0
M	35-44	>20	1 0	0 1 0 0 0	0 0 0 1	1 0 1 0 0 0 0 0 0 1
M	45-54	< 3	1 0	0 0 1 0 0	1 0 0 0	1 0 0 1 0 0 0 0 0 0
M	45-54	3-9	1 0	0 0 1 0 0	0 1 0 0	1 0 0 1 0 0 0 1 0 0
M	45-54	10-19	1 0	0 0 1 0 0	0 0 1 0	1 0 0 1 0 0 0 0 1 0
M	45-54	>20	1 0	0 0 1 0 0	0 0 0 1	1 0 0 1 0 0 0 0 0 1
M	55-64	< 3	1 0	0 0 0 1 0	1 0 0 0	1 0 0 0 1 0 0 0 0 0
M	55-64	3-9	1 0	0 0 0 1 0	0 1 0 0	1 0 0 0 1 0 1 0 0 0
M	55-64	10-19	1 0	0 0 0 1 0	0 0 1 0	1 0 0 0 1 0 0 0 1 0
M	55-64	>20	1 0	0 0 0 1 0	0 0 0 1	1 0 0 0 1 0 0 0 0 1
M	65-74	< 3	1 0	0 0 0 0 1	1 0 0 0	1 0 0 0 0 1 0 0 0 0
M	65-74	3-9	1 0	0 0 0 0 1	0 1 0 0	1 0 0 0 0 1 1 0 0 0
M	65-74	10-19	1 0	0 0 0 0 1	0 0 1 0	1 0 0 0 0 1 0 1 0 0
M	65-74	>20	1 0	0 0 0 0 1	0 0 0 1	1 0 0 0 0 1 0 0 1 0
F	25-34	< 3	0 1	1 0 0 0 0	1 0 0 0	1 1 0 0 0 0 0 0 0 0
F	25-34	3-9	0 1	1 0 0 0 0	0 1 0 0	1 1 0 0 0 0 0 1 0 0
F	25-34	10-19	0 1	1 0 0 0 0	0 0 1 0	1 1 0 0 0 0 0 0 1 0
F	25-34	>20	0 1	1 0 0 0 0	0 0 0 1	1 1 0 0 0 0 0 0 0 1
F	35-44	< 3	0 1	0 1 0 0 0	1 0 0 0	1 1 1 0 0 0 0 0 0 0
F	35-44	3-9	0 1	0 1 0 0 0	0 1 0 0	1 1 1 0 0 0 0 1 0 0
F	35-44	10-19	0 1	0 1 0 0 0	0 0 1 0	1 1 1 0 0 0 0 0 1 0
F	35-44	>20	0 1	0 1 0 0 0	0 0 0 1	1 1 1 0 0 0 0 0 0 1
F	45-54	< 3	0 1	0 0 1 0 0	1 0 0 0	1 1 0 1 0 0 0 0 0 0
F	45-54	3-9	0 1	0 0 1 0 0	0 1 0 0	1 1 0 1 0 0 0 1 0 0
F	45-54	10-19	0 1	0 0 1 0 0	0 0 1 0	1 1 0 1 0 0 0 0 1 0
F	45-54	>20	0 1	0 0 1 0 0	0 0 0 1	1 1 0 1 0 0 0 0 0 1
F	55-64	< 3	0 1	0 0 0 1 0	1 0 0 0	1 1 0 0 1 0 0 0 0 0
F	55-64	3-9	0 1	0 0 0 1 0	0 1 0 0	1 1 0 0 1 0 1 0 0 0
F	55-64	10-19	0 1	0 0 0 1 0	0 0 1 0	1 1 0 0 1 0 0 0 1 0
F	55-64	>20	0 1	0 0 0 1 0	0 0 0 1	1 1 0 0 1 0 0 0 0 1
F	65-74	< 3	0 1	0 0 0 0 1	1 0 0 0	1 1 0 0 0 1 0 0 0 0
F	65-74	3-9	0 1	0 0 0 0 1	0 1 0 0	1 1 0 0 0 1 1 0 0 0
F	65-74	10-19	0 1	0 0 0 0 1	0 0 1 0	1 1 0 0 0 1 0 1 0 0
F	65-74	>20	0 1	0 0 0 0 1	0 0 0 1	1 1 0 0 0 1 0 0 1 0
Chi-Square Statistic			1101.13	662.30	196.37	34.85
Degrees of Freedom			38	35	36	31
P-Value			< 0.001	< 0.001	< 0.001	.290

TABLE 8: LINEAR HYPOTHESES CONCERNING THE MODEL  $X_4$  FOR THE REGULAR DENTIST PERCENTAGES

HYPOTHESIS	CONTRAST MATRIX	CHI-SQUARE STATISTIC	D.F.	P-VALUE
No variation between sex subdomains	[0 1 0 0 0 0 0 0 0]	78.66	1	< 0.001
No variation among age subdomains	$\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$	75.69	4	< 0.001
No variation among income subdomains	$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$	530.61	3	< 0.001