

ANALYSIS OF CATEGORICAL DATA
OBTAINED BY STRATIFIED
RANDOM SAMPLING I.

by

E. Sobel, M. Francis and P. Imrey*

Department of Biostatistics
University of North Carolina at Chapel Hill

Institute of Statistics Mimeo Series No. 1065

March 1976

ANALYSIS OF CATEGORICAL DATA
OBTAINED BY STRATIFIED
RANDOM SAMPLING I.

E. Sobel, M. Francis and P. Imrey*

*E. Sobel is Research Fellow, Department of Epidemiology, and M. Francis is Assistant Professor, Department of Biostatistics, University of North Carolina, Chapel Hill, N. C. P. Imrey is Assistant Professor of Biostatistics, University of Illinois, Urbana, Ill.

ABSTRACT

A non-iterative method of analysis by linear models is extended to categorical data obtained by stratified random sampling. It is shown that, asymptotically, proportional allocation reduces the variance of estimators over that obtained by simple random sampling. The difference between the asymptotic covariance matrices of estimated parameter vectors obtained by simple random sampling and stratified random sampling with proportional allocation is shown to be positive definite under fairly non-restrictive conditions.

1. INTRODUCTION

Grizzle, Starmer and Koch (GSK, [4]) propose a weighted least-squares methodology for the analysis of a wide range of linear and non-linear models for categorical data, including testing of linear hypotheses within the underlying parameter space. Their approach is a synthesis of the work of Wald [11], Neyman [10] and Bhapkar [1,2]. The broad applicability of the method has been demonstrated in a continuing series of papers, e.g., Koch and Reinfurt [9], Forthofer and Koch [3], Koch, Johnson and Tolley [8], and Koch, Imrey and Reinfurt [7].

The research described above assumes throughout that the categorical data were obtained by independent simple random sampling within each of an arbitrary number of very large populations. Koch, Freeman and Freeman [6] discuss alternative sample designs for which this assumption may not be necessary due to the use of pseudo-replication or like methods for estimating the covariance matrix of the sample cell counts. Johnson and Koch [5], in the context of an example, show how proportions or mean scores from stratified random sampling designs may be analyzed within the original GSK context, through consideration of each stratum as a population and use of a "mediator matrix" to transform within-stratum cell counts or mean scores into the usual stratified estimates and their variances and covariances. Linear models are then applied to the statistics so derived.

In this paper we extend the suggestion of Johnson and Koch [5] to the full range of functional models described in GSK [4]. We then provide a multivariate extension of the well-known theoretic result that

stratified sampling with proportional allocation provides an unbiased estimator of a population mean with smaller variance than the estimator obtained from simple random sampling. In particular we then show that this result carries over to stratified and non-stratified estimators of the parameters in the functional models mentioned above.

2. NOTATION and the PROBLEM

As far as possible, we follow the notation used in GSK. We assume that there are s populations of interest. Within the i^{th} population we are interested in a multinomial random (response) variable with response categories C_{i1}, \dots, C_{ir_i} . Notice that we allow the number of categories to vary between populations. Let

π_{ij} = proportion of the i^{th} population
"belonging" to C_{ij}

$$\pi'_i = (\pi_{i1}, \dots, \pi_{ir_i})$$

$$\pi = (\pi'_1, \dots, \pi'_s)$$

The unknown π_{ij} are hypothesized to satisfy a possibly non-linear model of the form

$$\underset{u \times 1}{F}(\pi) = \underset{u \times v}{X} \underset{v \times 1}{\beta} \quad (1)$$

where

i) $\underset{u \times 1}{F}(\pi) = (f_1(\pi), \dots, f_u(\pi))$,

ii) $u \leq \sum (r_i - 1)$,

iii) each $f_m(\pi)$ has continuous second partials with respect to the π_{ij} ,

iv) the $f_m(\cdot)$ are independent of one another, i.e., the rank of

$$\underset{u \times \sum r_i}{H}(\pi) = \left(\frac{\partial f_m(\pi)}{\partial \pi_{ij}} \right)$$

is equal to u ,

v) $\underset{u \times v}{X}$ is a known $u \times v$ design matrix of rank $v \leq u$,

vi) $\underset{v \times 1}{\beta}$ is a $v \times 1$ vector of unknown parameters.

Assuming simple random sampling from each population, GSK treated (1) as a null hypothesis and presented an asymptotic χ^2 -test. Then, assuming (1) is valid, they gave a best asymptotically normal

(BAN) estimate of $\underline{\beta}$. Finally they gave an asymptotic χ^2 -test for the null hypothesis

$$\underline{C} \underline{\beta} = \underline{0} \quad , \quad (2)$$

where

vii) \underline{C} is a known $d \times v$ matrix of rank $d \leq v$.

A. NOTATION - Simple Random Sampling

Within the i^{th} population, $i = 1, \dots, s$, let

n_i = total sample size

n_{ij} = number of C_{ij} responses,
 $j = 1, \dots, r_i$ in the sample

$p_{ij} = n_{ij}/n_i$.

$\underline{p}'_i = (p_{i1}, \dots, p_{ir_i})$

$$\text{var}_{r_i \times r_i} (\underline{p}'_i) = \underline{V}(\underline{\pi}'_i) = \frac{1}{n_i} \begin{pmatrix} \pi_{i1}(1-\pi_{i1}) - \pi_{i1}\pi_{i2} & \dots & -\pi_{i1}\pi_{ir_i} \\ \vdots & & \vdots \\ -\pi_{ir_i}\pi_{i1} & -\pi_{ir_i}\pi_{i2} & \dots & \pi_{ir_i}(1-\pi_{ir_i}) \end{pmatrix}$$

Also let

$$\underset{1}{\overset{p}{\sim}} \underset{\Sigma r_i}{\times} = (p_1, \dots, p_s)$$

$$\underset{\Sigma r_i}{\overset{V(\pi)}{\sim}} \underset{\Sigma r_i}{\times} = \text{block diagonal matrix of the } \underset{\sim}{V}(\underset{\sim}{\pi}_i)$$

$$\underset{\sim}{V}(p) = \text{usual sample estimate of } \underset{\sim}{V}(\underset{\sim}{\pi})$$

$$\underset{\sim}{H} = \underset{\sim}{H}(p)$$

$$\underset{\sim}{S} = \underset{\sim}{H} \underset{\sim}{V}(p) \underset{\sim}{H}'$$

$$\underset{u}{\overset{S}{\sim}} \underset{x}{\times} \underset{u}{\sim} (\underset{\sim}{\pi}) = \underset{\sim}{H}(\underset{\sim}{\pi}) \underset{\sim}{V}(\underset{\sim}{\pi}) \underset{\sim}{H}(\underset{\sim}{\pi})'$$

B. NOTATION - Stratified Random Sampling

Within the i^{th} population, $i = 1, \dots, s$, let

$$n_{i.} = \text{total sample size}$$

$$D_i = \text{number of strata}$$

$$n_{i(h)} = \text{sample size in stratum } h$$

$$n_{i(h)j} = \text{number of } C_{ij} \text{ responses in the sample}$$

from stratum h , $j = 1, \dots, r_i$

$N_i.$ = total population size

N_{ih} = population size of stratum h

$W_{ih} = N_{ih}/N_i.$

$\tilde{W}_i' = (W_{i1}, \dots, W_{iD_i})$

$\alpha_{i(h)j}$ = proportion in stratum h
"belonging" to C_{ij}

$a_{i(h)j} = n_{i(h)j}/n_{i(h)}$

$\tilde{\alpha}'_{i(h)} = (\alpha_{i(h)1}, \dots, \alpha_{i(h)r_i}); \tilde{a}'_{i(h)} = (a_{i(h)1}, \dots, a_{i(h)r_i})$

$\tilde{\alpha}'_i = (\tilde{\alpha}'_{i(1)}, \dots, \tilde{\alpha}'_{i(D_i)}); \tilde{a}'_i = (\tilde{a}'_{i(1)}, \dots, \tilde{a}'_{i(D_i)})$
 $1 \times D_i r_i \quad 1 \times D_i r_i$

$\text{var}_{r_i \times r_i} (\tilde{a}_{i(h)}) = V(\tilde{\alpha}_{i(h)})$

$$= \frac{1}{n_{i(h)}} \begin{pmatrix} \alpha_{i(h)1}(1 - \alpha_{i(h)1}) - \alpha_{i(h)1} \alpha_{i(h)2} \dots - \alpha_{i(h)1} \alpha_{i(h)r_i} \\ - \alpha_{i(h)r_i} \alpha_{i(h)1} - \alpha_{i(h)r_i} \alpha_{i(h)2} \dots - \alpha_{i(h)r_i} (1 - \alpha_{i(h)r_i}) \end{pmatrix}$$

$\underset{\sim}{V}(\alpha_i)$ = block diagonal matrix of the $\underset{\sim}{V}(\alpha_{i(h)})$.
 $D_{i r_i}^{r_i} \times D_{i r_i}^{r_i}$

Also let

$$1 \times \underset{\sim}{\Sigma} D_{i r_i}^{\alpha'} = (\alpha'_1, \dots, \alpha'_s) \quad ; \quad 1 \times \underset{\sim}{\Sigma} D_{i r_i}^{\alpha'} = (\alpha'_1, \dots, \alpha'_s)$$

$\underset{\sim}{V}(\alpha)$ = block diagonal matrix of the $\underset{\sim}{V}(\alpha_i)$
 $\Sigma D_{i r_i}^{r_i} \times \Sigma D_{i r_i}^{r_i}$

$\underset{\sim}{V}(\hat{\alpha})$ = usual sample estimate of $\underset{\sim}{V}(\alpha)$

$$\Sigma_{i r_i}^{r_i} \times \underset{\sim}{\Sigma} D_{i r_i}^{\alpha'} = \begin{pmatrix} W_1^{-1} \otimes I_{r_1} & & & \\ & \ddots & & \\ & & \ddots & \\ & & & W_s^{-1} \otimes I_{r_s} \end{pmatrix}$$

where \otimes is the Kronecker product and I_{r_i} is the $r_i \times r_i$ identity matrix.

Finally let

$$S_{\sim st} = H_{\sim st} \underset{\sim}{V}(\hat{\alpha}) H_{\sim st}'$$

$$S_{\sim st}(\alpha) = H_{\sim st}(\alpha) \underset{\sim}{V}(\alpha) H_{\sim st}(\alpha)'$$

where $H_{\sim st}$ and $H_{\sim st}(\alpha)$ are defined in (3) in Section 4.

We remark here that even though $\underline{V}(\pi)$ and $\underline{V}(\alpha)$ are singular and thus only positive semi-definite, both $\underline{S}(\pi)$ and $\underline{S}_{st}(\alpha)$ are non-singular and thus positive definite as a consequence of assumption iv) following equation (1). We shall use the positive definiteness in the proof of the theorem in Section 6.

3. UNDER SIMPLE RANDOM SAMPLING

GSK pointed out that if (1) is valid, then

$$\underline{SS}[\underline{F}(\pi) = \underline{X}\beta] = (\underline{F}(p) - \underline{X}b)' \underline{S}^{-1} (\underline{F}(p) - \underline{X}b)$$

is asymptotically a $\chi^2(u-v)$ where

$$\underline{b} = (\underline{X}' \underline{S}^{-1} \underline{X})^{-1} \underline{X}' \underline{S}^{-1} \underline{F}(p).$$

This provides a test of (1). Note that if $u = v$ then (1) must trivially be true. If (1) is true the b is a BAN estimation of β . GSK suggested that to test (2) one use

$$\underline{SS}[\underline{C}\beta = \underline{0}] = \underline{b}' \underline{C}' (\underline{C}(\underline{X}' \underline{S}^{-1} \underline{X})^{-1} \underline{C}')^{-1} \underline{C} \underline{b}$$

which, if (2) is true, is asymptotically a $\chi^2(d)$. Note that if (1) is true and $\underline{F}(\pi)$ is linear, $\underline{S}(\pi)$ is the covariance matrix of $\underline{F}(p)$ and that if $\underline{F}(\pi)$ is non-linear

$$\underline{\text{var}} (\underline{F}(p)) \underline{S}^{-1}(\pi) \rightarrow \underline{I}$$

as the $n_i \rightarrow \infty$ so that $n_{i(h)}/n_i$ is constant. Similarly,

$$\tilde{S}\tilde{S}^{-1}(\tilde{\pi}) \rightarrow \tilde{I},$$

and $(\tilde{X}\tilde{S}^{-1}\tilde{X})^{-1}$ is asymptotically equal to the covariance matrix of \tilde{b} .

4. UNDER STRATIFIED RANDOM SAMPLING

In order to test (1), obtain a BAN estimator of β and test (2) under stratified random sampling we basically treat this situation from the point of view of simple random sampling from $\sum_{i=1}^S D_i$ populations.

Without loss of generality, we assume

$$\alpha_{i(h)j} > 0, \text{ for each } h, j, i.$$

Since

$$\pi_{ij} = \sum_{h=1}^{D_i} W_{ih} \alpha_{i(h)j}$$

we have

$$\tilde{Q}\tilde{\alpha} = \tilde{\pi}$$

and thus if (1) is true

$$\tilde{T}(\tilde{\alpha}) \equiv \tilde{F}(\tilde{Q}\tilde{\alpha}) = \tilde{X}\tilde{\beta}. \tag{1'}$$

Note that this differs from Johnson and Koch [5], who calculate functions within each stratum before pooling using stratum weights. Such an approach is inappropriate if $F(\cdot)$ is non-linear, while it produces results identical to ours in the linear case.

It is clear that, with $T(\cdot) \equiv F(Q\cdot)$, (1) and (1') are equivalent. Next we show that the conditions corresponding to iii) and iv) of section 2 are valid for

$$T(\alpha) = (t_1(\alpha), \dots, t_u(\alpha))'$$

We have

$$\frac{\partial t_m(\alpha)}{\partial \alpha_{i(h)j}} = \frac{\partial f_m(\alpha)}{\partial \alpha_{i(h)j}} = \frac{\partial f_m(\pi)}{\partial \pi_{ij}} \frac{\partial \pi_{ij}}{\partial \alpha_{i(h)j}} = W_{ih} \frac{\partial f_m(\pi)}{\partial \pi_{ij}}$$

and

$$\frac{\partial^2 t_m(\alpha)}{\partial \alpha_{i(h)j} \partial \alpha_{i'(h')j'}} = \frac{\partial [W_{ih} \partial f_m(\alpha) / \partial \pi_{ij}]}{\partial \alpha_{i'(h')j'}}$$

$$= W_{ih} W_{i'h'} \frac{\partial^2 f_m(\pi)}{\partial \pi_{ij} \partial \pi_{i'j'}}$$

Thus each $t_m(\alpha)$ has continuous second partials. Let

$$H_{\sim st \sim}(\alpha) = \left(\frac{\partial t_m(\alpha)}{\partial \alpha_{i(h)j}} \right)_{u \times \sum D_i r_i}$$

and

$$H_{\sim st} = H_{\sim st}(\alpha) = \left(\frac{\partial t_m(\alpha)}{\partial \alpha_{i(h)j}} \mid \alpha = \sim a \right)$$

Then

$$\begin{aligned} H_{\sim st}(\alpha) &= \left(\begin{array}{cccc} \left(\frac{\partial f_1(\pi)}{\partial \pi_1} \right)' W_{11} & \cdots & \left(\frac{\partial f_1(\pi)}{\partial \pi_1} \right)' W_{1D_1} & \cdots & \left(\frac{\partial f_1(\pi)}{\partial \pi_s} \right)' W_{s1} & \cdots & \left(\frac{\partial f_1(\pi)}{\partial \pi_s} \right)' W_{sD_s} \\ \vdots & & & & & & \vdots \\ \left(\frac{\partial f_u(\pi)}{\partial \pi_1} \right)' W_{11} & & \cdots & & & & \left(\frac{\partial f_u(\pi)}{\partial \pi_s} \right)' W_{sD_s} \end{array} \right) \\ &= \left(\begin{array}{ccc} W_{\sim 1} \otimes \left(\frac{\partial f_1(\pi)}{\partial \pi_1} \right)' & \cdots & W_{\sim s} \otimes \left(\frac{\partial f_1(\pi)}{\partial \pi_s} \right)' \\ \vdots & & \vdots \\ W_{\sim 1} \otimes \left(\frac{\partial f_u(\pi)}{\partial \pi_1} \right)' & \cdots & W_{\sim s} \otimes \left(\frac{\partial f_u(\pi)}{\partial \pi_s} \right)' \end{array} \right) \\ &= H(\pi)Q. \end{aligned}$$

Since \tilde{Q} has only one non-zero entry in each column and is $\sum r_i \times \sum D_i r_i$,

$$\text{rk } \tilde{Q} = \sum r_i.$$

Thus

$$\text{rk } H_{\tilde{st}}(\alpha) = u.$$

Consequently

$$\text{SS}[\tilde{T}(\alpha) = \tilde{X}\tilde{\beta}] = (\tilde{T}(\alpha) - \tilde{X}\tilde{b}_{\tilde{st}}) \tilde{S}_{\tilde{st}}^{-1} (\tilde{T}(\alpha) - \tilde{X}\tilde{b}_{\tilde{st}})$$

is asymptotically $\chi^2(u-v)$ where

$$\tilde{b}_{\tilde{st}} = (\tilde{X}' \tilde{S}_{\tilde{st}}^{-1} \tilde{X})^{-1} \tilde{X}' \tilde{S}_{\tilde{st}}^{-1} \tilde{T}(\alpha).$$

Thus if (1') is valid, $\tilde{b}_{\tilde{st}}$ is a BAN estimator of β . The test suggested by GSK for (2) then uses

$$\text{SS}[\tilde{C}\tilde{\beta} = \tilde{0}] = \tilde{b}_{\tilde{st}}' \tilde{C}' (\tilde{C}(\tilde{X}' \tilde{S}_{\tilde{st}}^{-1} \tilde{X})^{-1} \tilde{C}')^{-1} \tilde{C} \tilde{b}_{\tilde{st}}$$

which is asymptotically $\chi^2(d)$ if (2) is true.

5. COMPARISON OF \tilde{b} and $\tilde{b}_{\tilde{st}}$ with PROPORTIONAL ALLOCATION

We now assume that

$$n_{i(h)} = W_{ih} n_{i.}, \quad h = 1, \dots, D_i, \quad i = 1, \dots, s.$$

Then the usual estimate of π_{ij} is

$$\hat{\pi}_{ij} = \frac{D_i}{\sum_{h=1}^{D_i} W_{ih}} a_{i(h)j}$$

Put $\hat{\pi}'_i = (\hat{\pi}_{i1}, \dots, \hat{\pi}_{ir_i})$ and $\hat{\pi}' = (\hat{\pi}'_1, \dots, \hat{\pi}'_s)$. Then

$$\underline{p}_{st} = \hat{\pi} = \underline{Q}\underline{a}$$

and

$$\begin{aligned} \underline{\text{var}}(\hat{\pi}) &= \underline{\text{var}}(\underline{p}_{st}) = (\text{cov}(\hat{\pi}_{ij}, \hat{\pi}_{i'j'})) \\ &= \underline{Q}\underline{V}(\underline{a})\underline{Q}' \end{aligned}$$

Also, in this case, since

$$\underline{H}_{st}(\underline{a}) = \underline{H}(\underline{\pi})\underline{Q},$$

$$\underline{S}_{st}(\underline{a}) = \underline{H}(\underline{\pi})\underline{\text{var}}(\hat{\pi})\underline{H}(\underline{\pi})'$$

and

$$\underline{S}_{st} = \underline{H}\underline{V}(\underline{p}_{st})\underline{H}'$$

where

$$\underline{V}(\underline{p}_{st}) = \text{usual estimate of } \underline{\text{var}}(\underline{p}_{st}).$$

As a notational convenience, $\underline{A} \leq \underline{B}$ and $\underline{A} < \underline{B}$, where \underline{A} and \underline{B} are square matrices, means that $\underline{B} - \underline{A}$ is positive semi-definite and positive definite, respectively. \geq and $>$ have equivalent meanings.

Lemma

$$\underline{\text{var}}(\hat{\underline{\pi}}) \leq \underline{V}(\underline{\pi}).$$

Proof:

First note that both $\underline{\text{var}}(\hat{\underline{\pi}})$ and $\underline{V}(\underline{\pi})$ are composed of (square) block diagonal matrices:

$$\underline{V}(\underline{\pi}) = \begin{pmatrix} \underline{V}(\underline{\pi}_1) & & \underline{0} \\ & \ddots & \\ \underline{0} & & \underline{V}(\underline{\pi}_s) \end{pmatrix}$$

and

$$\underline{\text{var}}(\hat{\underline{\pi}}) = \begin{pmatrix} \underline{\text{var}}(\hat{\underline{\pi}}_1) & & \underline{0} \\ & \ddots & \\ \underline{0} & & \underline{\text{var}}(\hat{\underline{\pi}}_s) \end{pmatrix}$$

Thus if we show that

$$\underline{\text{var}}(\hat{\underline{\pi}}_i) \leq \underline{V}(\underline{\pi}_i), \quad i = 1, \dots, s$$

we are done. But

$$n_i \cdot (\underline{V}(\underline{\pi}_i) - \underline{\text{var}}(\hat{\underline{\pi}}_i)) =$$

$$\begin{pmatrix} \sum_h w_{ih} (\alpha_{i(h)1} - \pi_{i1})^2 & \sum_h w_{ih} (\alpha_{i(h)1} - \pi_{i1})(\alpha_{i(h)2} - \pi_{i2}) \dots \sum_h w_{ih} (\alpha_{i(h)1} - \pi_{i1})(\alpha_{i(h)r_i} - \pi_{ir_i}) \\ \vdots & \vdots \\ \sum_h w_{ih} (\alpha_{i(h)r_i} - \pi_{ir_i})(\alpha_{i(h)1} - \pi_{i1}) & \sum_h w_{ih} (\alpha_{i(h)r_i} - \pi_{ir_i})(\alpha_{i(h)2} - \pi_{i2}) \dots \sum_h w_{ih} (\alpha_{i(h)r_i} - \pi_{ir_i})^2 \end{pmatrix}$$

Now let

$$\underline{X}'_i = (X_{i1}, \dots, X_{ir_i})$$

$$X_{ij} = \alpha_{i(h)j}, \quad j = 1, \dots, r_i$$

where h is a randomly chosen integer, $1 \leq h \leq D_i$, with

$$P(h = h_0) = W_{ih_0}$$

Then

$$E(X_{ij}) = \sum_h W_{ih} \alpha_{i(h)j} = \pi_{ij}$$

$$\text{var}(X_{ij}) = \sum_h W_{ih} (\alpha_{i(h)j} - \pi_{ij})^2$$

and

$$\text{cov}(X_{ij}, X_{ij'}) = \sum_h W_{ih} (\alpha_{i(h)j} - \pi_{ij})(\alpha_{i(h)j'} - \pi_{ij'})$$

Thus the covariance matrix of \underline{X}_i is $n_i \cdot (\underline{V}(\underline{\pi}_i) - \underline{\text{var}}(\hat{\underline{\pi}}_i))$, which is therefore positive semi-definite.

Theorem

Asymptotically, with proportional allocation,

$$\text{var } \underset{\sim}{\underset{\sim}{c}} \hat{b}_{st} \leq \text{var } \underset{\sim}{\underset{\sim}{c}} \hat{b}$$

for all vectors $\underset{\sim}{c} \neq \underset{\sim}{0}$. This is equivalent to

$$(\underset{\sim}{X}' \underset{\sim}{S}_{st}^{-1} \underset{\sim}{X})^{-1} \leq (\underset{\sim}{X}' \underset{\sim}{S}^{-1} \underset{\sim}{X})^{-1}.$$

Proof:

Since

$$\text{var } \underset{\sim}{\underset{\sim}{c}} \hat{b}_{st} = \underset{\sim}{c}' (\underset{\sim}{X}' \underset{\sim}{S}_{st}^{-1} \underset{\sim}{X})^{-1} \underset{\sim}{c}$$

$$\text{var } \underset{\sim}{\underset{\sim}{c}} \hat{b} = \underset{\sim}{c}' (\underset{\sim}{X}' \underset{\sim}{S}^{-1} \underset{\sim}{X})^{-1} \underset{\sim}{c}$$

the equivalence is clear. Now from the lemma and the fact that $\underset{\sim}{H}(\pi)$ has linearly independent rows it follows that

$$\underset{\sim}{S}(\pi) - \underset{\sim}{S}_{st}(\alpha) = \underset{\sim}{H}(\pi) [\underset{\sim}{V}(\pi) - \underset{\sim}{\text{var}}(\hat{\pi})] \underset{\sim}{H}'(\pi) \geq \underset{\sim}{0}.$$

Since $\underset{\sim}{S}_{st}(\alpha) > \underset{\sim}{0}$, its inverse possesses a (symmetric) square root $\underset{\sim}{P}_1(\alpha) = \underset{\sim}{P}_1$, say for convenience. Then $\underset{\sim}{P}_1 \underset{\sim}{S}(\pi) \underset{\sim}{P}_1' > \underset{\sim}{0}$ and by the Principal Axis Theorem there exists an orthogonal matrix $\underset{\sim}{P}_2(\alpha) = \underset{\sim}{P}_2$, say, such that

$$\underset{\sim}{P}_2 (\underset{\sim}{P}_1 \underset{\sim}{S}(\pi) \underset{\sim}{P}_1') \underset{\sim}{P}_2' = \underset{\sim}{D},$$

where $\underset{\sim}{D}$ is diagonal of full rank. Thus, with $\underset{\sim}{P}(\alpha) = \underset{\sim}{P} = \underset{\sim}{P}_2 \underset{\sim}{P}_1$,

$$\underset{\sim}{P} \underset{\sim}{S}(\pi) \underset{\sim}{P}' = \underset{\sim}{D}$$

$$\underset{\sim}{P} \underset{\sim}{S}_{st}(\alpha) \underset{\sim}{P}' = \underset{\sim}{P}_2 \underset{\sim}{P}_2' = \underset{\sim}{I}.$$

But $\underset{\sim}{D} - \underset{\sim}{I} = \underset{\sim}{P} [\underset{\sim}{S}(\pi) - \underset{\sim}{S}_{st}(\alpha)] \underset{\sim}{P}' \geq \underset{\sim}{0}$ and therefore every diagonal element of $\underset{\sim}{D}$ is greater than or equal to 1. Let

$$\underset{\sim}{E} = \underset{\sim}{I} - \underset{\sim}{D}^{-1}.$$

Then

$$\begin{aligned} (\underline{X}' \underline{S}_{st}^{-1} (\underline{\alpha}) \underline{X})^{-1} &= (\underline{X}' \underline{P}' \underline{P} \underline{X})^{-1} = [(\underline{P} \underline{X})' (\underline{P} \underline{X})]^{-1} \\ (\underline{X}' \underline{S}^{-1} (\underline{\pi}) \underline{X})^{-1} &= (\underline{X}' \underline{P}' \underline{D}^{-1} \underline{P} \underline{X})^{-1} = [(\underline{P} \underline{X})' (\underline{I} - \underline{E}) (\underline{P} \underline{X})]^{-1} \end{aligned}$$

Put

$$\underline{U}(\underline{\delta}) = [(\underline{P} \underline{X})' (\underline{I} - \underline{E}(\underline{\delta})) (\underline{P} \underline{X})]^{-1}$$

where

$$\underline{\delta} \in \underline{\Sigma} = \left\{ \underline{\delta} = (\delta_1, \dots, \delta_u) \mid 0 \leq \delta_k < 1, k = 1, \dots, u \right\}$$

$\underline{E}(\underline{\delta})$ = diagonal matrix with $\underline{\delta}$ as the diagonal vector

Thus

$$\begin{aligned} \underline{U}(0) &= (\underline{X}' \underline{S}_{st}^{-1} \underline{X})^{-1} \\ \underline{U}(\underline{\delta}_{\pi}) &= (\underline{X}' \underline{S}^{-1} (\underline{\pi}) \underline{X})^{-1} \end{aligned}$$

for some $\underline{\delta}_{\pi} \in \underline{\Sigma}$. We are then finished once we prove that $\underline{x}' \underline{U}(\underline{\delta}) \underline{x}$, $\underline{x} \neq 0$, is a non-decreasing function of each δ_k for $\underline{\delta} \in \underline{\Sigma}$. But

$$\begin{aligned} \frac{\partial [\underline{x}' \underline{U}(\underline{\delta}) \underline{x}]}{\partial \delta_k} &= \frac{\underline{x}' \left[\frac{\partial \underline{U}(\underline{\delta})}{\partial \delta_k} \right] \underline{x}}{\left[\frac{\partial \underline{U}(\underline{\delta})}{\partial \delta_k} \right]} \\ &= -\underline{x}' \underline{U}(\underline{\delta}) (\underline{P} \underline{X})' \left[\frac{\partial (\underline{I} - \underline{E}(\underline{\delta}))}{\partial \delta_k} \right] (\underline{P} \underline{X}) \underline{U}(\underline{\delta}) \underline{x} \end{aligned}$$

$$\begin{aligned}
 &= [\underline{u}(\delta)_{\underline{x}}] \underline{\theta}_k \underline{\theta}_k' [\underline{u}(\delta)_{\underline{x}}] \\
 &= [\underline{\theta}_k' \underline{u}(\delta)_{\underline{x}}] [\underline{\theta}_k' \underline{u}(\delta)_{\underline{x}}] \\
 &\geq 0.
 \end{aligned}$$

where $\underline{\theta}_k$ is the k^{th} column of $\underline{P}\underline{X}$. This completes the proof.

Corollary

Let $\underline{X}' = (\underline{X}'_1, \dots, \underline{X}'_s)$ where \underline{X}_i is defined in the proof of the lemma.

Let \underline{X} be a matrix whose $\prod_{i=1}^s D_i$ rows are the realizations of \underline{X} .

Then $\text{var } \underline{c}' \underline{b}_{st} < \text{var } \underline{c}' \underline{b}$ for all $\underline{c} \neq \underline{0}$ if and only if

$$\text{rank} [\underline{\ell}, \underline{X}\underline{H}'] = u + 1,$$

where $\underline{\ell}$ is a column vector of 1's.

Proof:

Since $\underline{S}(\underline{\pi}) - \underline{S}_{st}(\underline{\alpha})$ is the covariance matrix of $\underline{H}\underline{X}$ under the probability mass function

$$P(\underline{X}' = (\underline{\alpha}'_1(h_1), \dots, \underline{\alpha}'_s(h_s))) = \prod_1^s P(\underline{X}_i = \underline{\alpha}_i(h_i)) = \prod_1^s W_i(h_i)$$

the condition given is equivalent to $\underline{S}(\underline{\pi}) - \underline{S}_{st}(\underline{\alpha}) > \underline{0}$. The corollary

then follows directly from the proof of the theorem.

6. DISCUSSION

Under the model $\underline{F}(\underline{\pi}) = \underline{X}\underline{\beta}$, where $\underline{F}(\cdot)$ satisfies the conditions given in Section 2, the estimator of $\underline{\beta}$ under the least-squares algorithm behaves asymptotically as a linear function of the estimators of the cell probabilities (see [10]). Consequently if $\underline{F}(\cdot)$ is such that its components, $f_k(\cdot)$, are derived from non-overlapping groups of the s populations, then the results in Section 5 are rather obvious. We have shown that this extends to rather arbitrary linear models.

The approach taken above ignores the finite population correction factor. However, as Johnson and Koch [5] remark, multiplication of $\underline{V}(p)$ by $(1-f)$, where f is the sampling fraction, results in valid large-sample statistics based on the work of Wald [11] when the population is considered finite. If the fpc were considered here the basic results would clearly remain valid when all covariance matrices are multiplied by the appropriate $(1-f_i)$. This is easily seen by consideration of sequences of problems with stratum sizes N_{ih} and sample sizes $n_{i(h)}$, with both N_{ih} and $n_{i(h)}$ increasing such that their ratio approaches $f_i = n_i/N_i$. The resulting hypergeometric distributions obey the same central limit theory as their multinomial counterparts.

The weighted least-squares procedure is vulnerable to small sample problems due to observed zero cell counts. In the unstratified case if this happens S will not be of full rank because the corresponding π_{ij} will be estimated by 0. GSK [4] suggest using an "estimate" $1/r_i n_i$. However, in the stratified case the invertibility of \underline{S}_{st} depends upon the covariance matrix of the stratified estimates of the π_{ij} , not upon

the covariance matrix from the individual strata, i.e., the covariance matrix of the estimates of the $\alpha_{i(h)j}$ for fixed i, h . Thus, problems will be encountered when all strata within a given population have the same cell empty. Following GSK, each corresponding $\alpha_{i(h)j}$ can be estimated by $1/r_i n_{i(h)}$. Such an occurrence presumably would arise no more often than having a zero cell count using simple random sampling with identical total sample size.

Finally, the formulation and results in Sections 4 and 5 allow for a general study of optimal allocation as a function of costs associated with sampling from each stratum, the model $\underline{F}(\underline{\pi}) = \underline{X}\underline{\beta}$ and the variance-covariance matrix of the estimators. Results along these lines will be presented in a sequel.

REFERENCES

- [1] Bhapkar, V.P., "Notes on Analysis of Categorical Data," Series No. 477, Institute of Statistics, University of North Carolina, 1966 (mimeographed).
- [2] Bhapkar, V.P., "A Note on the Equivalence of Two Test Criteria for Hypotheses in Categorical Data," Journal of the American Statistical Association, 61 (March 1966) 228-35.
- [3] Forthofer, Ronald N. and Koch, Gary G., "An Analysis for Compounded Functions of Categorical Data," Biometrics, 29 (March 1973), 143-57.
- [4] Grizzle, James, E., Starmer, C. Frank and Koch, Gary G., "Analysis of Categorical Data by Linear Models," Biometrics, 25 (September 1969), 489-503.
- [5] Johnson, William D. and Koch, Gary G., "Analysis of Qualitative Data: Linear Models," Health Sciences Research, (Winter 1970), 358-69.
- [6] Koch, Gary G., Freeman, Daniel H. Jr., and Freeman, Jean L., "Strategies in the Multivariate Analysis of Data from Complex Surveys," International Statistical Review, 43 (No. 1, 1975), 59-78.
- [7] Koch, Gary G., Imrey, Peter B., and Reinfurt, Donald W., "Linear Model Analysis of Categorical Data with Incomplete Response Vectors," Biometrics, 28 (September 1972), 663-92.
- [8] Koch, Gary G., Johnson, William D., and Tolley, H. Dennis, "A Linear Models Approach to the Analysis of Survival and Extent of Disease in Multidimensional Contingency Tables," Journal of the American Statistical Association, 67 (December 1972), 783-96.
- [9] Koch, Gary G. and Reinfurt, Donald W., "The Analysis of Categorical Data from Mixed Models," Biometrics, 27 (March 1971), 157-73.
- [10] Neyman, Jenzy, "Contribution to the Theory of the χ^2 Test," in J. Neyman, Ed., Proceedings of the Berkeley Symposium on Mathematical Statistics and Probability, University of California Press, Berkeley and Los Angeles, 1949, 239-73.
- [11] Wald, Abraham, "Tests of Statistical Hypotheses Concerning Several Parameters where the Number of Observations Is Large," Transactions of the American Mathematics Society, 54 (November 1943), 426-482.