

Leonard A. Marascuilo, University of California at Berkeley

Maryellen McSweeney, Michigan State University

## Introduction

Multistage sampling procedures and nonresponse of sampled units frequently make the analysis of data generated by analytical surveys extremely difficult. If there are only two independent subpopulations of particular interest, the analysis of survey data is not especially complex since the subpopulation parameters may be estimated from the data. If in addition, the sample sizes are large, it is always possible to test the null hypothesis,

$$H_0: \theta_1 = \theta_2$$

against the alternative,

$$H_1: \theta_1 \neq \theta_2$$

by computing

$$Z = \frac{\hat{\theta}_1 - \hat{\theta}_2}{\sqrt{SE_{\hat{\theta}_1}^2 + SE_{\hat{\theta}_2}^2}}$$

and referring the observed Z value to standard normal curve tables. If  $H_0$  is rejected, a point estimate of the parameter difference is given by  $(\hat{\theta}_1 - \hat{\theta}_2)$  and a  $(1 - \alpha)\%$  interval estimate is given by

$$(\hat{\theta}_1 - \hat{\theta}_2) \pm Z_{\alpha/2} \sqrt{SE_{\hat{\theta}_1}^2 + SE_{\hat{\theta}_2}^2}$$

For studies involving more than two subpopulations comparable analytical methods have not been reported. Investigation of the technical literature shows that Gold (1963) and Goodman (1964) have extended the simultaneous confidence interval method of Scheffé (1959) to certain special cases associated with the parameters of contingency tables and the parameters of Markov Chains. Marascuilo (1966) has extended their model to include multiple confidence intervals for correlation coefficients and for sample averages from analysis of variance designs in which the variances are unequal. Since analogous situations occur in survey research studies, this extension should be of considerable value in the analysis of survey data.

In this paper a proof of the chi-square analog of Scheffé's Theorem is given. From the results of this proof a simple-to-compute test statistic is proposed for the test of the hypothesis

$$H_0: \theta_1 = \theta_2 = \dots = \theta_K = \theta_0$$

against the alternative

$$H_1: H_0 \text{ is false.}$$

The post hoc multiple comparison procedures associated with the rejection of the hypothesis  $H_0$  are indicated. These methods are used to test and identify the sources of differences in attitudes expressed by adult citizens toward the integration of de facto segregated schools in three different socio-economic subpopulations of an urban American community.

## Chi-square Analog of Scheffé's Theorem

Consider a univariate model in which there are K treatments, conditions, or populations. Let the parameters of the model be represented by  $\theta' = (\theta_1, \theta_2, \dots, \theta_K)$ . Let  $\hat{\theta}' = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_K)$  be a set of large sample efficient estimators of these unknown parameters. Furthermore, let the covariance matrix for these estimators be of rank  $q \leq K$ . It is known from large sample theory that

$$U = (\hat{\theta} - \theta)' (\text{Cov}(\hat{\theta}))^{-1} (\hat{\theta} - \theta)$$

has an asymptotic chi-square distribution with q degrees of freedom since it is the exponent in the asymptotic K-variate normal distribution of  $\hat{\theta}' = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_K)$ . Since the rank of U is q, it is possible to find a set of q linearly independent estimable functions

$$\psi_p = \sum_{k=1}^K a_k \theta_k = \underline{a}' \underline{\theta}$$

that will span the space of all contrasts of the form

$$\psi = \sum_{p=1}^q c_p \psi_p = \underline{c}' \underline{\psi} \quad p = 1, 2, \dots, q$$

Let the set of all possible contrasts of the form  $\psi$  be denoted by L. A set of estimates for the  $\psi_p$  is

$$\hat{\psi}_p = a_1 \hat{\theta}_1 + a_2 \hat{\theta}_2 + \dots + a_K \hat{\theta}_K = \underline{a}' \hat{\underline{\theta}}$$

Since the  $\hat{\psi}_p$  are linearly independent functions of asymptotically multi-variate normal random variables with a covariance matrix of rank q, they are also asymptotically multivariate normal with exponent given by

$$Q = (\hat{\underline{\psi}} - \underline{\psi})' (\text{Cov}(\hat{\underline{\psi}}))^{-1} (\hat{\underline{\psi}} - \underline{\psi})$$

As a result  $Q$  must have an asymptotic chi-square distribution with  $q$  degrees of freedom. Thus, a  $(1 - \alpha)\%$  confidence ellipsoid for the point  $\underline{\psi}' = (\psi_1, \psi_2, \dots, \psi_q)$  is given by  $Q \leq \chi_q^2 (1 - \alpha)$ .

This confidence ellipsoid serves as the basis for the analog of Scheffé's Theorem. The proof of this theorem parallels, as one would expect, the proof of Scheffé's Theorem. The notation used is that of Scheffé so that the two proofs may be easily compared. The analog of Scheffé's Theorem reads as follows:

**Theorem.** The probability is  $(1 - \alpha)$  in the limit that simultaneously for all  $\underline{\psi} \in L$

$$\hat{\underline{\psi}} - \sqrt{\chi_q^2 (1 - \alpha)} \sqrt{\text{Var}(\hat{\underline{\psi}})} \leq \underline{\psi} \leq \hat{\underline{\psi}} + \sqrt{\chi_q^2 (1 - \alpha)} \sqrt{\text{Var}(\hat{\underline{\psi}})}$$

**Proof.** The inequality that defines the asymptotic confidence ellipsoid for the point

$$\underline{\psi}' = (\psi_1, \psi_2, \dots, \psi_q) \text{ is}$$

$$(\hat{\underline{\psi}} - \underline{\psi})' (\text{Cov}(\hat{\underline{\psi}}))^{-1} (\hat{\underline{\psi}} - \underline{\psi}) \leq \chi_q^2 (1 - \alpha)$$

The point  $\underline{\psi}$  is in the ellipsoid if, and only if, it lies between all pairs of parallel planes of support of the ellipsoid. If  $\underline{c}' = (c_1, c_2, \dots, c_q)$  is an arbitrary nonzero vector, Scheffé has shown that the point  $\underline{\psi}$  lies between the two planes of support of the ellipsoid orthogonal to  $\underline{c}$  if, and only if,

$$|\underline{c}' (\underline{\psi} - \hat{\underline{\psi}})| \leq \sqrt{\underline{c}' \underline{M}^{-1} \underline{c}}$$

In this case,

$$\underline{M} = \frac{1}{\chi_q^2 (1 - \alpha)} (\text{Cov}(\hat{\underline{\psi}}))^{-1}$$

Thus

$$\begin{aligned} \underline{c}' \underline{M}^{-1} \underline{c} &= \underline{c}' \left( \frac{1}{\chi_q^2 (1 - \alpha)} (\text{Cov}(\hat{\underline{\psi}}))^{-1} \right)^{-1} \underline{c} \\ &= \underline{c}' \chi_q^2 (1 - \alpha) \text{Cov}(\hat{\underline{\psi}}) \underline{c} \\ &= \chi_q^2 (1 - \alpha) [\underline{c}' \text{Cov}(\hat{\underline{\psi}}) \underline{c}] \end{aligned}$$

Since any contrast  $\underline{\psi}$  in  $L$  can be estimated by

$$\hat{\underline{\psi}} = \sum_{p=1}^q c_p \hat{\underline{\psi}}_p = \underline{c}' \hat{\underline{\psi}},$$

the variance of the estimate is given by

$$\text{Var}(\hat{\underline{\psi}}) = \underline{c}' \text{Cov}(\hat{\underline{\psi}}) \underline{c}.$$

Therefore

$$\underline{c}' \underline{M}^{-1} \underline{c} = \chi_q^2 (1 - \alpha) \text{Var}(\hat{\underline{\psi}})$$

and

$$\sqrt{\underline{c}' \underline{M}^{-1} \underline{c}} = \sqrt{\chi_q^2 (1 - \alpha)} \sqrt{\text{Var}(\hat{\underline{\psi}})}$$

Since

$$|\underline{c}' \underline{\psi} - \underline{c}' \hat{\underline{\psi}}| = |\underline{\psi} - \hat{\underline{\psi}}|$$

the last inequality actually states that simultaneously for all  $\underline{\psi} \in L$  the probability in the limit is  $(1 - \alpha)$  that

$$|\underline{\psi} - \hat{\underline{\psi}}| \leq \sqrt{\chi_q^2 (1 - \alpha)} \sqrt{\text{Var}(\hat{\underline{\psi}})}$$

This completes the proof.

As with the F-test, the Chi-square test will reject  $H_0$  if, and only if, the estimate  $\hat{\underline{\psi}}$  of at least one  $\underline{\psi}$  is significantly different from zero. Equivalently, if  $H_0$  is rejected there is at least one contrast in the  $\hat{\underline{\psi}}_p$  that is significantly different from zero.

#### Derivation of the test statistic to test

$$H_0: \theta_1 = \theta_2 = \dots = \theta_K = \theta_0$$

In most surveys the subpopulations consist of population strata or domains of investigation that usually comprise mutually exclusive subsets of the total universe of interest. For this reason, the estimates of the parameters within the individual subpopulations are statistically independent, so that  $\text{Cov}(\hat{\theta}_i, \hat{\theta}_j) = 0$  for  $i \neq j$  and

$$U = (\hat{\underline{\theta}} - \underline{\theta})' (\text{Cov}(\hat{\underline{\theta}}))^{-1} (\hat{\underline{\theta}} - \underline{\theta})$$

reduces to

$$\sum_{k=1}^K \frac{(\hat{\theta}_k - \theta_k)^2}{\text{Var}(\hat{\theta}_k)}$$

which is asymptotic  $\chi_K^2$  with  $q = K$ .

To test the hypothesis

$$H_0: \theta_1 = \theta_2 = \dots = \theta_K = \theta_0$$

it is only necessary to evaluate  $U$  under  $H_0$  and determine whether or not  $U > \chi_K^2 (1 - \alpha)$ . If  $U$  is too large,  $H_0$  is rejected.

For most applications, the exact value of  $\theta_0$  is unknown and must be estimated. An easy-to-obtain estimate is the one that minimizes  $U$ . This estimate is given by

$$\hat{\theta}_o = \frac{\sum_{k=1}^K \frac{1}{\text{Var}(\hat{\theta}_k)} \hat{\theta}_k}{\sum_{k=1}^K \frac{1}{\text{Var}(\hat{\theta}_k)}} = \frac{\sum_{k=1}^K w_k \hat{\theta}_k}{\sum_{k=1}^K w_k}$$

If this estimate is substituted into U and if the resulting expression is denoted as  $U_o$ , it follows that

$$U_o = \sum_{k=1}^K \frac{(\hat{\theta}_k - \hat{\theta}_o)^2}{\text{Var}(\hat{\theta}_k)} = \sum_{k=1}^K w_k (\hat{\theta}_k - \hat{\theta}_o)^2$$

Familiar analysis of variance methods can be employed to show that  $U_o$  is asymptotically chi-square with  $(K - 1)$  degrees of freedom. Therefore, a simple decision rule that may be used for testing  $H_o$  is: reject  $H_o$  if  $U_o > \chi^2_{K-1}(1 - \alpha)$  and do not reject  $H_o$  if  $U_o < \chi^2_{K-1}(1 - \alpha)$ .

If the variances are unknown and the sample sizes are large, the large sample estimates of the variances can be substituted into the final result with little loss. This also applies to the estimate of  $\hat{\theta}_o$  which would then be equal to

$$\hat{\theta}_o = \frac{\sum_{k=1}^K \frac{1}{\text{var}(\hat{\theta}_k)} \hat{\theta}_k}{\sum_{k=1}^K \frac{1}{\text{var}(\hat{\theta}_k)}} = \frac{\sum_{k=1}^K \hat{w}_k \hat{\theta}_k}{\sum_{k=1}^K \hat{w}_k}$$

In addition, the test statistic would be

$$U_o' = \sum_{k=1}^K \hat{w}_k (\hat{\theta}_k - \hat{\theta}_o)^2$$

### Examples

By means of a principal component analysis based on 1960 census data, the 28 census tracts of Berkeley, California, were partitioned into three mutually exclusive subpopulations representing high, medium, and low socio-economic status areas. Within each census tract a two per cent sample of adults was selected. The following sampling procedure was used to obtain independent samples for each census tract. A city block was chosen at random with probability proportional to the block size reported in the 1960 census data, and a simple random sample of six households was

taken for the selected block. This process was repeated for additional blocks within a census tract until the number of adults in the sample was estimated to equal two per cent of the 1960 census tract adult population. The sampling procedure was repeated in each of the 28 census tracts. Since the population magnitude varied across census tracts, the sample sizes over census tracts ranged from 25 to 86 adults. The removal of wrong addresses and vacant houses from the sampling frame reduced the actual sampling fraction to 1.88 per cent and the initial sample size to 1,392 adults.

The survey was begun on April 15, 1964, with the mailing of letters and questionnaires to the 742 randomly selected households in the community. On April 29, follow-up letters were mailed to all nonresponding households. Between May 11th, and 19th, a random sample of 1/3 of the remaining nonrespondents was interviewed by trained female personnel from the Survey Research Center of the University of California. Usable information was ultimately obtained from 971 adults of the originally selected sample.

One of the items appearing on the questionnaire read as follows:

For some (elementary) schools the committee suggested that lines be changed so that the percentage of nonwhite and white children in these schools would be more like the percentage for the entire school system.

- (1) \_\_\_ I agree
- (2) \_\_\_ I disagree
- (3) \_\_\_ I am not sure

The "I am not sure" category of response was excluded in the analysis of the data. The analysis of this item considered the effect of socio-economic status on attitudes toward increasing the racial integration of the schools by means of boundary changes. In particular, it was hypothesized that members of the low SES Negro areas of the community would show the strongest support for the boundary changes designed to effect school integration while the greatest opposition would be expressed by the high SES white areas.

The sampling unit for this survey was the household, but the unit for analysis was the individual respondent. Consequently the number of adults per household was a random variable, the value of which was undetermined until data were obtained for each household. Since the number of adults per household was unknown prior to sampling, the per cent agreeing to the change in school boundaries was estimated by a separate ratio estimate,  $\hat{P}_h$ , for each census tract. Furthermore separate ratio estimates,  $\hat{P}_{hi}$ , were required for each wave of response within a single census tract because the responses to the original letter, the follow-up letter, or the personal interview produced an artificial stratification of the respondents for each census tract. Despite the small sample sizes

within strata, no appreciable differences between the separate and combined ratio estimates were found. Separate ratio estimates were chosen in preference to combined estimates on the basis of greater simplicity of computation and explication. As a result, the final parameter estimate for each of the three subpopulations defined by principal component analysis involved primary stratification of the census tracts together with the artificial within-tract stratification based on the wave of response.

For a subpopulation defined by principal component analysis:

$$1. \hat{p} = \sum_{h=1}^L \left( \frac{n_h}{n} \right) \hat{p}_h \quad h = 1, 2, \dots, L \text{ census tracts}$$

where

$$2. \hat{p}_h = \left( \frac{n_{h1}}{n_h} \right) \hat{p}_{h1} + \left( \frac{n_{h2}}{n_h} \right) \hat{p}_{h2} + \left( \frac{n_h - n_{h1} - n_{h2}}{n_h} \right) \hat{p}_{h3} \quad i = 1, 2, 3 \text{ waves of response}$$

and

$$3. \hat{p}_{hi} = \frac{\sum_{j=1}^{n_{hi}} a_{hij}}{\sum_{j=1}^{n_{hi}} m_{hij}} \quad \begin{array}{l} j = 1, 2, \dots, n_{hi} \\ \text{households in census tract } h \text{ that} \\ \text{answer in the } i\text{th.} \\ \text{wave of response} \\ a_{hij} = \text{number of} \\ \text{adults in household } j \\ \text{who answered "I agree."} \\ m_{hij} = \text{number of} \\ \text{adults in household } j \end{array}$$

The approximate variance is estimated by

$$4. SE_{\hat{p}}^2 = \sum_{h=1}^L \left( \frac{n_h}{n} \right)^2 SE_{\hat{p}_h}^2$$

where

$$5. SE_{\hat{p}_h}^2 = \left( \frac{n_{h1}}{n_h} \right)^2 SE_{\hat{p}_{h1}}^2 + \left( \frac{n_{h2}}{n_h} \right)^2 SE_{\hat{p}_{h2}}^2 + \left( \frac{n_h - n_{h1} - n_{h2}}{n_h} \right)^2 SE_{\hat{p}_{h3}}^2$$

and

$$6. SE_{\hat{p}_{hi}}^2 = \frac{1}{n_{hi} \bar{m}_{hi}^2} \sum_{j=1}^{n_{hi}} \frac{(a_{hij} - \hat{p}_{hi} m_{hij})^2}{n_{hi} - 1}$$

In Table 1 the distribution of response to the question by subpopulation is shown. If binomial estimates of the variances are used, the hypothesis of equal proportions agreeing in the three subpopulations will be rejected since  $\chi^2 = 94.80$  exceeds  $\chi^2_{2(.95)} = 5.99$ . However, binomial estimates and the chi-square test of homogeneity are inappropriate because the responses within a cluster (household) are not independent, but positively correlated.

The appropriate ratio estimates of the parameters for the three subpopulations are given in Table 2. For each subpopulation the estimated variance of the proportion agreeing is considerably larger for the ratio estimate than for the corresponding binomial estimate.

For these data,

$$\hat{p}_0 = \frac{\sum_{k=1}^3 \hat{w}_k \hat{p}_k}{\sum_{k=1}^3 \hat{w}_k} \quad k = 1, 2, 3 \text{ subpopulations defined by principal component analysis}$$

$$\hat{p}_0 =$$

$$\frac{691.1039(.862) + 394.4510(.586) + 652.4490(.348)}{691.1039 + 394.4510 + 652.4490}$$

$$.606$$

and

$$U_0' = 691.1039(.862 - .606)^2 + 394.4510(.586 - .606)^2 + 652.4490(.348 - .606)^2$$

$$U_0' = 88.88$$

Table 1: Distribution of Responses by Subpopulation and Binomial Estimates of the Parameters Based on Proportional Allocation

<u>Response</u>	<u>Subpopulation</u>			<u>Total</u>
	<u>Low SES</u>	<u>Medium SES</u>	<u>High SES</u>	
Agree	170	179	105	454
Disagree	30	126	157	313
Total	200	305	262	767
Per Cent Agreement (binomial estimate)	.850	.586	.401	.592
Variance (binomial estimate)	.00063	.00079	.00091	

Table 2: Ratio Estimates of the Parameters by Subpopulation

<u>Parameter Estimated</u>	<u>Subpopulation</u>			<u>Total</u>
	<u>Low SES</u>	<u>Medium SES</u>	<u>High SES</u>	
Per Cent Agreement	.862	.586	.348	.606
Variance	.001447	.002535	.001533	
Weight, $\hat{w}_k$	691.1039	394.4510	652.4490	

Where the estimate of the total is given by  $\hat{p}_0 = \frac{\sum_{k=1}^K \hat{w}_k \hat{p}_k}{\sum_{k=1}^K \hat{w}_k} = .606$

Table 3: 95% Confidence Intervals for the Set of Simple Contrasts

<u>Contrast</u>	<u>Value of Contrast <math>\hat{p}_k - \hat{p}_{k'}</math></u>	<u>Estimated Variance of Contrast</u>	<u>Lower Limit of Confidence Interval</u>	<u>Upper Limit of Confidence Interval</u>	<u>Signifi- cance</u>
Low vs. Medium	.862-.586	.001447 + .002535	.122	.430	Sig.
Low vs. High	.862-.348	.001447 + .001533	.377	.651	Sig.
Medium vs. High	.586-.348	.002535 + .001533	.082	.394	Sig.

Since  $U'_0 > \chi^2_{2(.95)} = 5.99$ ,  $H_0$  is rejected. Thus there is reason to believe that at least one linear contrast of the parameters is significantly different from zero.

For this study, the general form of the  $\begin{pmatrix} 3 \\ 2 \end{pmatrix}$  or 3 simple contrasts is given by

$$\hat{\psi}_{kk'} = \hat{p}_k - \hat{p}_{k'}$$

$$k \neq k'$$

with the estimated variance given by

$$\text{var}(\hat{\psi}_{kk'}) = \text{var}(\hat{p}_k) + \text{var}(\hat{p}_{k'})$$

These contrasts and their estimated variances are summarized in Table 3. All three contrasts are statistically significant from zero at the overall .05 level.

Although these hypothesis testing and multiple contrast techniques have been illustrated for the case of three independent subpopulations, their range of possible application in analytical surveys is far broader. For example, the hypothesis of equality of a set of domain means could be tested by these techniques. If the domains are defined by the strata of a stratified sampling procedure, the estimates of the domain means and of their variances given by Cochran (1963, pp. 148-149) could be substituted into the test statistic  $U'_0$ . If the hypothesis of equal domain means is rejected because  $U'_0 > \chi^2_{K-1}(1 - \alpha)$ , then statistically significant sources of differences could be determined by use of the post hoc procedure suggested in this paper.

Furthermore it should be noted that the general theorem permits one to test hypotheses and determine simultaneous confidence intervals for analytical surveys in which the parameter estimates are not independent. An example of correlated ratio estimates in a survey in which the sampling unit consists of clusters of households is suggested by Cochran (1963, p. 182). A test of the hypothesis that the proportion of men who smoke is equal to the proportion of women who smoke could be based on  $U'_0$ . The test statistic would be given by

$$U'_0 = (\hat{\theta}_1 - \hat{\theta}_0, \hat{\theta}_2 - \hat{\theta}_0) \begin{pmatrix} \text{var}(\hat{\theta}_1) & \text{cov}(\hat{\theta}_1, \hat{\theta}_2) \\ \text{cov}(\hat{\theta}_1, \hat{\theta}_2) & \text{var}(\hat{\theta}_2) \end{pmatrix}^{-1} \begin{pmatrix} \hat{\theta}_1 - \hat{\theta}_0 \\ \hat{\theta}_2 - \hat{\theta}_0 \end{pmatrix}$$

where the estimate of  $\theta_0$  which minimizes  $U'_0$  would be

$$\hat{\theta}_0 = \frac{\hat{\theta}_1 \text{var}(\hat{\theta}_2) + \hat{\theta}_2 \text{var}(\hat{\theta}_1) - (\hat{\theta}_1 + \hat{\theta}_2) \text{cov}(\hat{\theta}_1, \hat{\theta}_2)}{\text{var}(\hat{\theta}_1) + \text{var}(\hat{\theta}_2) - 2 \text{cov}(\hat{\theta}_1, \hat{\theta}_2)}$$

The extension of this test to three or more domains could also be based on  $U'_0$  where the elements of the covariance matrix could be obtained by the formulas given by Keyfitz (1957) or by Kish and Hess (1959). If the hypothesis of equality of the ratios were to be rejected, sources of the differences in the parameters could be determined by the post hoc procedure outlined above. The estimated variances of linear contrasts in the ratios could be obtained by substitution of the elements of the covariance matrix into the formula

$$\text{var}(\hat{\psi}) = \sum_{k=1}^K \text{var}(\hat{\theta}_k) + \sum_{k \neq k'} \text{cov}(\hat{\theta}_k, \hat{\theta}_{k'})$$

## Summary

The analysis of data generated by analytical surveys is compounded by complex sampling procedures and the nonresponse of sampled units. The problem is significantly greater when the number of subpopulations of interest exceeds two. On the basis of a chi-square analog of Scheffé's Theorem a simple multiple contrast or confidence interval procedure can be generated that can be used to identify possible parameter differences provided that the null hypothesis of no parameter differences has been rejected. This method should prove to be of considerable use to scientists whose major research methodology involves survey sampling.

## References

- COCHRAN, WILLIAM G. (1963). Sampling Techniques. New York: John Wiley and Sons, Inc.
- GOOD, R. Z. (1963). Tests auxiliary to  $\chi^2$  tests in a Markov Chain. Ann. Math. Statist. 31, 56-74.
- GOODMAN, LEO A. (1964). Interactions in multidimensional contingency tables. Ann. Math. Statist. 35, 716-725.
- KEYFITZ, NATHAN (1957). Estimates of sampling variance where two units are selected from each stratum. J. Amer. Statist. Ass. 52, 503-510.
- KISH, LESLIE & HESS, IRENE (1959). On variances of ratios and their differences in multistage samples. J. Amer. Statist. Ass. 54, 416-446.
- MARASCUILO, LEONARD A. (1966). Large-sample multiple comparisons. Psychol. Bull. 65, 280-290.
- SCHEFFÉ, HENRY (1959). The Analysis of Variance. New York: John Wiley and Sons, Inc.

## Footnotes

1. This report was prepared at the Institute of Human Learning which is supported by grants from the National Science Foundation and the National Institutes of Health.