Barry L. Ford, U.S. Department of Agriculture R. R. Hocking, Mississippi State University Anne Coleman, Mississippi State University

PURPOSE:

The purpose of this paper is to discuss certain sampling plans in the context of respondent burden and their application to a large scale agricultural survey. These survey procedures attempt to lower the respondent burden by utilizing multivariate regression.

INTRODUCTION:

Surveys sometimes use questionnaires which require a great deal of time for the respondent to complete because of a large number of questions. This situtation occurs not only because the subject may be complex but also because the additional cost of asking 100 questions, for example, instead of 10 questions is relatively small. Once a survey is designed and implemented, budget constraints often dictate obtaining as much information as possible from each respondent.

For example, the U.S. Department of Agriculture (USDA) runs a survey to estimate the expenditures of farmers when producing their crops and livestock. The questionnaire for this survey contains over 600 questions--questions ranging from expenditures on seeds to expenditures for diesel oil to expenditures for wire. The amount of detail needed on farm expenditures is so great that the ideal solution is to use a set of independent surveys, each of which obtains information on a different part of a farmer's expenditures. However, the number of farmers contacted would then be enourmous, and thus, the cost of getting expenditure information would also be enourmous. Therefore, the USDA carries out one survey with a questionnaire which asks about all expenditures but may require five or six hours to complete.

Recently the burden placed on respondents by surveys has received more attention because of complaints by the public about the large number of surveys (both in government and in private industry) and the inconvenience of these surveys. These complaints cause the USDA to be more aware of the respondent burden caused by its surveys. Respondent burden now has no exact definition in the literature but is often used in a general sense to signify the time and other costs to each respondent answering a survey questionnaire. When the USDA obtains expenditure information by avoiding a set of surveys with a small number of questions on each survey and by using one survey with a large number of questions, the USDA changes a small respondent burden on many people to a large respondent burden on a few people. One could also view this trade-off as an exchange of monetary cost to the government for an increase in respondent burden. Since everyone benefits from a reduction in government costs, everyone benefits from this exchange except that sample of farmers who must give the time and effort to answer over 600 questions.

Can some survey procedure reduce the burden placed on this sample of farmers without greatly

increasing costs to the survey organization? Under certain conditions the answer is yes. The usual circumstance is that the questions on a survey all relate to one specific area (eg. expenditures of farmers) so that the variables obtained from each respondent are highly correlated with each other. Therefore, the correlation structure of the variables can be used to avoid asking <u>all</u> of the respondents <u>every</u> question and still achieve a satisfactory level on the estimates.

SAMPLE DESIGNS:

Suppose the goal is to make estimates on a large set of variables which are denoted by the vector w. If w is composed of two sets of

variables,
$$\underline{w} = \begin{bmatrix} \underline{x} \\ \underline{y} \end{bmatrix}$$
, where $\underline{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_p \end{bmatrix}$ and $\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_q \end{bmatrix}$,

then let plan A be the use of two independent samples to measure \underline{x} and \underline{y} . If the size of each sample size is n, then the data would be of the form:

1	-
<u>x</u> 1	
$\frac{x_1}{2}$	
$\frac{x}{2}$	1
-	-
:	1
•	ì
$\frac{\mathbf{x}}{\mathbf{n}}$	
	_!
v	ł
y _{n+1}	
y _{n+2}	
-1172	
	į
•	i
Σ _{2n}	
	_

where the subscript refers to data collected on a particular observation (i.e. respondent).

To minimize survey cost an organization might obtain all of the variables on each respondent. One can call this strategy plan B:

$\frac{x_1}{x_2}$	<u>у</u> 1 У ₂
:	:
<u>x</u> n	y _n

A simple method of using the correlation structure to avoid asking all n respondents every question is to use a multivariate regression estimator in a double sampling approach. The variables in \underline{x} should be correlated with the variables composing y, and thus \underline{x} can fairly well predict <u>y</u> by using a multivariate linear regression. This sample design, plan C, would obtain values for both <u>x</u> and <u>y</u> for the first n_1 observa-

tions and then obtain values of only \underline{x} on the remaining $(n-n_1)$ observations:

$$\begin{array}{c|c} \underline{x}_{1} & \underline{y}_{1} \\ \underline{x}_{2} & \underline{y}_{2} \\ \vdots \\ \underline{x}_{n_{1}} & \underline{y}_{n_{1}} \\ \hline \underline{x}_{n_{1}+1} & \underline{x}_{n_{1}+2} \\ \hline \underline{x}_{n_{1}+2} & \vdots \\ \vdots \\ \underline{x}_{n} & \underline{x}_{n} \end{array}$$

To estimate the mean vector of \underline{w} , $\underline{\mu}_{w}$; one estimates the two components- $\underline{\mu}_{w} = \begin{bmatrix} \underline{\mu}_{x} \\ \underline{\mu}_{y} \end{bmatrix}$. Of course, the estimator of $\underline{\mu}_{x}$ is the average of n vectors, i.e.:

$$\hat{\underline{\mu}}_{\mathbf{X}} = \frac{\hat{\underline{\lambda}} \times \hat{\underline{\lambda}}_{\mathbf{i}}}{n}$$

To estimate $\underline{\mu}_y$, however, one must apply a multivariate linear regression estimator: If $\underline{\times}_x$ and \underline{y}_y are averages using only the first n respondents, then:

$$\hat{\underline{\mu}}_{y} = \underline{\underline{y}}^{*} + B^{*} (\hat{\underline{\mu}}_{x} - \underline{\underline{x}}^{*})$$

where B is a matrix of regression coefficient given by B = $\sum_{xx}^{-1} \sum_{xy}$ and:

 $\Sigma = \begin{bmatrix} \Sigma_{\mathbf{x}\mathbf{x}} & \Sigma_{\mathbf{x}\mathbf{y}} \\ \Sigma_{\mathbf{x}\mathbf{y}}^{\prime} & \Sigma_{\mathbf{y}\mathbf{y}} \end{bmatrix}$

is the variance - covariance matrix of $\underline{w}.$ The variance of $\underline{\mu}_{\mathbf{v}}$ is, of course, :

$$\operatorname{Var}(\hat{\underline{\mu}}_{X}) = \frac{1}{n} \Sigma_{XX}.$$

Given a known B, the variance of μ_v is [1]:

$$\operatorname{Var}(\hat{\underline{\mu}}_{y}) = \frac{1}{n} \Sigma_{yy} - \frac{(n-n_{1})}{nn_{1}} \quad B^{*} \Sigma_{xx} B.$$

Obviously, plan C will cause the standard errors of the variables in <u>y</u> to be larger than those in plan A or B. A re-expression [1] of Var (μ_v) is:

Var
$$(\hat{\mu}_{i}) = \begin{bmatrix} 1 & -\frac{(n-n_{1})}{n} & R_{i}^{2} \end{bmatrix} \sigma_{i}^{2}/n_{1}$$

for $i = 1, 2, ...q$

where μ_i is the estimate of $E(y_i)$. One can easily observe that the change in the level of the variance of the estimates is closely related to R_i , the multiple correlation coefficient for variable i. Thus, if <u>x</u> is a fairly good predicator of <u>y</u>, there is little loss in accuracy. The R_i^2 (for i = 1, 2, ..., q) should determine

whether one should adopt plan C over plan B. What is the respondent burden in plans A, B, and C? There are several different ways one might wish to measure respondent burden, but for the purposes of this paper we simply adopt the number of questions per respondent. If B represents respondent burden, one has:

$$B_{A} = \frac{\frac{1}{2} (np + nq)}{n}$$
$$B_{B} = \frac{np + nq}{n}$$
$$B_{C} = \frac{np + nq}{n}$$

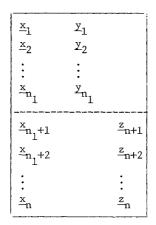
The enumerators of the burdens are written so that obviously $\mathcal{B}_A < \mathcal{B}_B$ if p and q are non-zero. Of course, \mathcal{B}_A would be even smaller and the cost of plan A even larger if more than two samples were used. Although plan C would not differ much in cost from plan B, its respondent burden is less than plan B. However, plan C does suffer a loss in the accuracy of estimating μ_y -- a loss which depends on R_i (i = 1, 2, ..., q). Because n_1 can be smaller when the R_i 's are larger, the respondent burden can be decreased when the R_i are larger.

Given a large number of variables, one might wish to elaborate on this double sampling scheme of plan C by separating the dependent variables into several closely related groups. For example, in the USDA expenditure survey there might be two groups of dependent variables:

group I : dependent variables relating to seed and plant expenses.

group II: dependent variables relating to fertilizers and pesticides.

The variables in \underline{x} should be highly related to the variables in both groups I and II. For instance, some of the variables in \underline{x} might be crop acreages. Groups I and II can be obtained from different respondents in the sample. Rather than just a <u>y</u> vector, there is a <u>y</u> vector and a <u>z</u> vector, and the sample design (called plan D) is:



Plan D is still a double sampling scheme, but there are now two sets of dependent variables.

One has $\underline{w} = \begin{bmatrix} \underline{x} \\ \underline{y} \\ \underline{z} \end{bmatrix}$ and $\underline{\mu}_{w} = \begin{bmatrix} \underline{\mu}_{x} \\ \underline{\mu}_{y} \\ \underline{\mu}_{z} \end{bmatrix}$

Also, Σ , the variance - covariance matrix of \underline{w} ; is now:

	Σ _{xx}	Σ_{xy}	Σ_{xz}
Σ =	Σ_ xy	Σуу	Σ_{yz}
	Σ'_{xz}	Σyz	Σzz

The estimators are of the same form as those in plan C:

$$\hat{\underline{\mu}}_{\mathbf{x}} = \hat{\underline{\Sigma}}_{\mathbf{x}} \underline{x}_{\mathbf{i}}/\mathbf{n}$$

$$\hat{\underline{\mu}}_{\mathbf{y}} = \hat{\underline{y}} + \mathbf{B}_{\mathbf{y}} (\hat{\underline{\mu}}_{\mathbf{x}} - \hat{\underline{x}}_{\mathbf{y}})$$

$$\hat{\underline{\mu}}_{\mathbf{z}} = \hat{\underline{z}} + \mathbf{B}_{\mathbf{z}} (\hat{\underline{\mu}}_{\mathbf{x}} - \hat{\underline{x}}_{\mathbf{z}})$$

where:

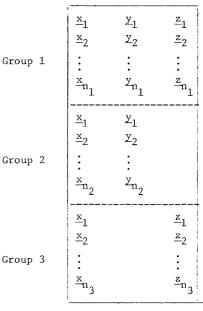
The variances of $\hat{\underline{\mu}}_{x}$, $\hat{\underline{\mu}}_{y}$ and $\hat{\underline{\mu}}_{z}$ are also of the same form as in plan C.

If $n_2 = (n-n_1)$, <u>x</u> has p variables, <u>y</u> has q variables, and <u>z</u> has r variables, then:

$$B_{\rm D} = \frac{n_{\rm I} (p+q) + n_{\rm 2} (p+r)}{n} = \frac{n_{\rm P} + n_{\rm 1} q + n_{\rm 2} r}{n} \,.$$

(When $n_1 = n_2$, then $\mathcal{B}_D = \mathcal{B}_C$.) The respondent burden in plan D does not differ greatly from plan C. However, the respondent burden in plan D is spread more evenly throughout the sample. Plan C has n_1 respondents answering all the questions and the rest of the sample answering a few questions while plan D has all respondents answering about the same number of questions.

Studying plan D, one might consider whether correlations exist between y and z which might be useful. If they do exist, one might want to ask a very small subgroup in the sample all of the questions in order to estimate Σ and use the information in Σ_{yz} . Thus, one has plan E:



To discuss the estimators for plan E, one must discuss a general class of sample designs studied by Hocking, Hartley, et al. [1, 2, 3]. The USDA calls sample designs in this class pattern samples because they depend on a pattern of missing values for certain observations. The observations in a pattern sample can be divided into groups based upon which variables are observed and which are not. These groups are called pattern groups. The first pattern group in plan E is a complete group---all variables are observed. The second group has only x and y observed, and third group has only \underline{x} and \underline{z} observed. A pattern design is called nested when the pattern groups can be arranged so that the variables in each pattern group are a subset of the preceding group. Plan C is a nested design, but plans D and E are not.

Each pattern group can be represented by a design matrix. A design matirx is composed of 0's and 1's to signify which variables are being observed. For example, if pattern group 1 is composed of four variables, its design matrix, D_1 , is the identity matrix:

	1	0	0	0
D. =	0	1	0	0
D ₁ =	0 0 0	0	1	0 0 0 1
	0	0	0	1
	1			

If the second pattern group only contains variables 2 and 4, then:

$$D_2 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

Pattern groups may contain variables which are totals of other variables. For instance, the third group might contain variable 1 and the sum of variables 3 and 4:

 $D_3 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}$

(Note that the rank of a design matrix can not be less than the number of rows.)

If the data are separated into pattern groups, then in the tth group there are n_t observations which yield an estimate of a mean vector, $\hat{\mu}_t$, of the variables in that group and $\hat{\Sigma}_t$, the estimated variance - covariance matrix among those variables. Under the assumption that the first group is a complete group, i.e. $D_1 = I$,

the maximum likelihood estimators are [1]:

where:

$$\hat{\mu}_{t} = D_{t} \hat{\mu}$$

 $\dot{\tilde{\Sigma}}_{t}$ = the sub-matrix of $\dot{\tilde{\Sigma}}$ pertaining to the variables in the tth group

 $H_{t} = (\hat{\underline{\mu}}_{t} - \hat{\underline{\mu}}_{t})(\hat{\underline{\mu}}_{t} - \hat{\underline{\mu}}_{t})^{2}$ $M_{t} = n_{t}(\hat{\hat{z}}_{t} + H_{t}).$

Equations I and II must be solved iteratively:

- 1: calculate $\underline{\mu}^{\times}$ by using the variance covariance matrix from the first pattern group as $\underline{\Sigma}^{\times}$.
- 2: use the $\frac{\star}{\mu}$ from step 1 to calculate H_t and M_t in order to form a new $\frac{\star}{\Sigma}$.
- 3: return to step 1 and use the new estimate of $\overset{*}{\Sigma}$ to make a new estimate $\overset{*}{\text{of }}\mu$.
- 4: keep cycling through steps 1-3 until a tolerance level is reached.

The maximum likelihood equations actually have a more general form which is not dependent on $D_1 = I$, but because of the additional notation

involved, this form is not discussed in this paper. The computer programs which perform the iterations (written by Anne Coleman [1, Attachment 2]) rely on $D_1 = I$.

Covergence in the iteration process is not assured in general, but in practical applications convergence has taken less than ten iterations.

The variance - covariance matrix of $\hat{\underline{\mu}}$ is estimated by:

$$W = \sum_{t=1}^{T} D_{t} W_{t} D_{t}$$

where

$$W_t = n_t \sum_{t=1}^{t-1}$$

Since there is no explicit formulation for the standard error of a particular variable, probably, the best method of seeing whether gains are likely in using plan E over plan D is to compare the multiple correlation coefficients. For the variables in y, for example, one can compare the multiple correlation coefficient when the variables in \underline{x} are used as independent variables to the multiple correlation coefficient when \underline{x} and \underline{z} are used as independent variables. Although not exactly correct, this approach should be a good indication of the improvement one might expect in adopting plan E over plan D.

The respondent burden of plan E is a slight increase of plan D if the percentage of respondents in the first pattern group is small. The burden is:

$$B_{E} = \frac{n_{1} (p+q+r) + n_{2} (p+q) + n_{3} (p+r)}{n}$$
$$= \frac{np + (n_{1}+n_{2}) q + (n_{1}+n_{3}) r}{n} \cdot$$

APPLICATION:

The purpose of this section is to compare a measure of respondent burden and the coefficients of variation resulting from sampling plans A, B, C, D, and E by using data from an agricultural survey. A complete data set from this survey was obtained under sampling plan B, but the values of certain variables in different sample units were simulated as missing in order to produce plans D and plan E.

The agricultural survey used as an illustration is a survey which estimates the various costs to farmers of growing agricultural products. This survey uses an elaborate questionnaire to obtain values for over 600 variables. This paper will actually restrict itself to a subset of 60 variables because of the computer cost involved. Although the agricultural survey is a complex sample design in practice, for illustrative purposes the data set is used as though it were from a simple random sample. Further work is planned in extending results to more complex sample designs.

Table 1 compares the respondent burden and the average coefficient of variation in plans A, B, C, D, and E. The measure of the respondent burden is the average number of variables obtained from each respondent, i.e. sample unit.

Plan A is obviously the optimum plan except that it costs approximately twice as much as the other plans. Thus, plan B is better than plan A, but the respondent burden is double. Plan C is better than plan B because there is only a small increase in the average coefficient of variation but the respondent is reduced by 35%. Plan D is better than plan C because it spreads the respondent burden more evenly across the entire sample. Plan E does decrease the average coefficient of variation slightly, but at a slight increase in both respondent burden and its spread. When one also consideres the more complicated estimation process involved with plan E, plan D is the most attractive sampling scheme.

Table 1: A comparison of five sampling plans in relation to data from an agricultural survey.

Plan	Respondent Burden	Spread of the Respondent Burden (Standard Deviation)	Average Coefficient of Variation	Approx. Cost
А	30	0	0.17	2(\$k)
В	60	0	0.17	\$k
С	39	21	0.21	\$k
D	39	1	0.21	\$k
E	41	7	0.20	\$k

BIBLIOGRAPHY:

- Hocking, R. R. Final Report on Design of Sample Surveys to Reduce Respondent Burden.
 1977. (A report written by contract with Research Division Economics, Statistics, and Cooperatives Service, Room 4818, South Building, U.S. Department of Agriculture, Washington, D.C. 20250.)
- [2] Hartley, H. O. and Hocking, R. R. "The Analysis of Incomplete Data," Biometrics, Volume 27, pages 783 - 824. 1971.
- [3] Hocking, R. R., Huddleston, H. F., and Hunt, H. H. "A Procedure for Editing Survey Data," Journal of the Royal Statistical Society -Series C, Volume 23, No. 2 pages 121 - 133. 1974.