

A SIMULATION STUDY TO COMPARE PROCEDURES WHICH IMPUTE FOR MISSING ITEMS ON AN ESS HOG SURVEY

BY

BARRY L. FORD, DOUGLAS G. KLEWENO, ROBERT D. TORTORA

**STATISTICAL RESEARCH DIVISION
ECONOMICS AND STATISTICS SERVICE
U.S. DEPARTMENT OF AGRICULTURE
WASHINGTON, D. C. 20250
DECEMBER 1980
ESS Staff Report No. AGE8801222.1**

A Simulation Study to Compare Procedures Which Impute for Missing Items on an ESS Hog Survey. By Barry L. Ford, Douglas G. Kleweno, and Robert D. Tortora; Statistical Research Division; Economics and Statistics Service; U.S. Department of Agriculture; Washington D.C. 20250; December, 1980. ESS Staff Report No. AGE8801222.1

ABSTRACT

This simulation study compares the effects of six procedures which impute values for missing items. An ESS list survey to estimate hog variables is used in this study. Specifically, the effects of the procedures on two survey variables--first quarter expected farrowings and second quarter expected farrowings--are analyzed under a range of conditions. These conditions include the percentage of missing items and the designation of which values are missing. Analysis of the mean square errors, the effects on the correlations, and the cost show that two versions of a ratio give the test results for very large sample sizes.

This paper was reproduced for limited distribution to the research community outside the U.S. Department of Agriculture. The views expressed herein are not necessarily those of ESS or USDA.

ACKNOWLEDGEMENT

The assistance of Ann Adams in programming one of the imputation procedures is gratefully acknowledged.

CONTENTS

	Page
Introduction	1
The Procedures	2
The Ratio Procedure	2
The Array Procedure	3
The Estmat Procedure	4
The Zero Spike Procedure	5
The IPrincomp Procedure ..	6
Analysis	6
Experimental Design	7
Results	8
Summary	15
Recommendations	16
Bibliography	16

1. INTRODUCTION

The problem of incomplete data is one of the most common problems of survey work. Incomplete data is of two types -- missing units and missing items. Missing units are the result of total nonresponse for a sample unit and, thus, consist of refusals and inaccessibilities. Missing items refer to those units which have missing values but also have some reported values. For example, the respondent answers some questions but not others, or he answers some questions incorrectly so that the true values are unknown. *The purpose of this study is to compare six procedures which impute for missing items and which could serve as alternatives to the operational procedure of hand editing by field statisticians.*

The problem of missing units is the subject of previous studies by the Economics and Statistics Service (ESS), of the U.S. Department of Agriculture [2]. The problem of missing items however, has received little attention in ESS-Statistics. The extent of field editing for missing items had never been quantified until a recent tabulation by Methods Staff of the 1979 December Enumerative Survey (DES) and of the December 1979 Multiple Frame Hog Survey. In the ten major hog producing states over fifteen percent of the sample with a nonzero number of hogs had at least one value imputed. First and second quarter expected farrowings, the two survey variables involved in this study, were imputed by field statisticians in over thirteen percent of the hog reports with a nonzero inventory. The 1979 DES had imputed values to the hog questions on over thirteen percent of the reports. Thus, the problem of missing items is large enough to warrant a research study.

The basic research tool of this study is simulation. Using a complete data set (no missing values) from the list sample of a multiple frame hog survey, the authors simulate which values are missing. Six missing item procedures are then applied, and the imputed values are compared to each other and to the original values. Although simulation experiments are in a sense artificial, they do allow analysis over a wide range of conditions and a comparison against "true" values.

The simulations in this study are over various levels of two effects:

- 1) the randomization mechanism used to designate which values are missing and
- 2) the rate at which values are missing.

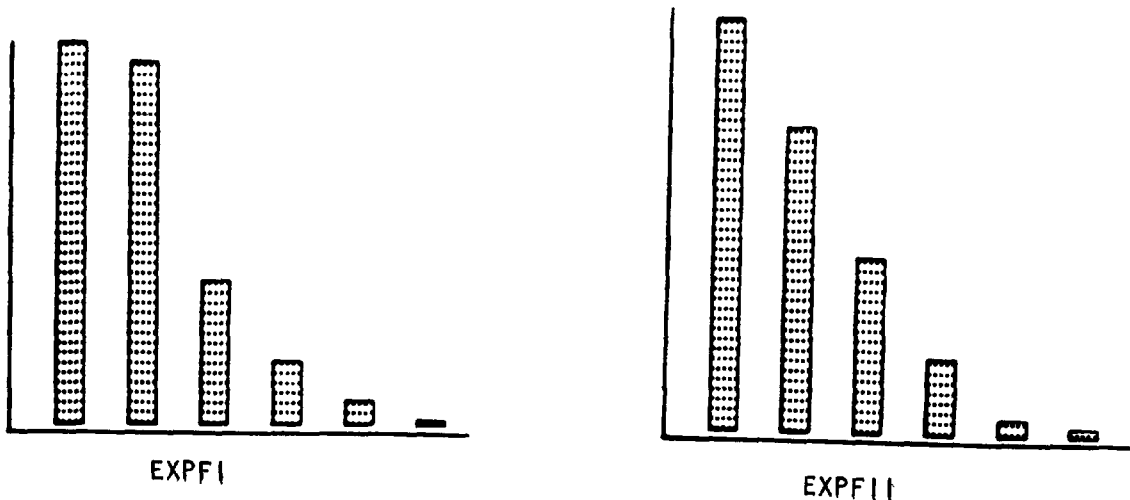
For each level of these two effects, there are several incomplete data sets simulated from the original data set. The original data set is divided into three replicates, and this replicate structure is carried over into each simulated data set. Missing item procedures are applied to each replicate independently in order to obtain unbiased estimates of standard errors.

The data set used in this study is from a March 1978 list sample in Iowa, a major hog producing state. Specifically, hog data from stratum 5 with boundary values of 600 to 999 hogs is used. Stratum 5 contains 201 complete sample units which are divided into three replicates of 67 units each. Each unit has 15 quantitative variables--sows, boars, other pigs, market hogs (five weight groups), total inventory, expected farrowings first quarter (EXPF1), expected farrowings second quarter (EXPF2), sows farrowed during previous three months, pigs now on hand, pigs sold, and one control variable used for stratifying the population. For the purposes of this study the authors confine the simulation of missing values to two hog variables, EXPF1 and EXPF2. Values

for either or both of these expected farrowing variables can be designated as missing. Also, all imputations must obey the edit check that the sum of EXPFI and EXPFII is less than or equal to total sows.

The distributions from the original data set for the EXPFI and EXPFII variables are both highly skewed. Figure 1.1 gives two bar graphs to show the general shape of the distributions. The first bar of each graph shows the number of sampling units with a zero value. The frequency of zeros is not a surprise due to the nature of the data being collected. Pork producers are not expected to have farrowings every quarter of the year. The mean of the first quarters expected farrowings is 22.11 and the variance 509.32; the mean of the second quarters expected farrowings is 21.45 and the variance 502.22. Thus, EXPFI and EXPFII are quite similar in distribution. The correlation between total sows and EXPFI is 0.82 and between total sows and EXPFII is 0.81. Of course, both EXPFI and EXPFII have integer values greater than or equal to zero.

Figure 1.1: Bar graphs to show the general shape of the distributions of EXPFI and EXPFII.



2. THE PROCEDURES

This section gives a description of the six procedures which impute for missing items. This description includes the estimation techniques and assumptions of the procedures. The descriptions are written in general terms of how the procedures would impute for a hypothetical data set which has both complete and incomplete units.

2.1 The Ratio Procedure (Variations 1 and 2)

The ratio procedure examined in this study imputes a value for each missing value by using the equation:

$$y_{\text{ratio}} = \hat{R} x^*$$

where:

\hat{R} is the estimated ratio between the variables x and y

x^* is the value of an x variable for a sample unit which has a missing y value

y_{ratio} is the value imputed for the missing y value.

An estimate of R is based on the sample units which are complete. If x' and y'

are totals for the complete units in the sample, then $\hat{R} = \frac{\sum y}{\sum x}$. Thus, this estimator of R assumes that the ratio for the complete units is a good estimate of the ratio for the incomplete units.

Although called an auxiliary variable, the x variable may be a survey variable or the control variable. When a y value is missing, the ratio procedure uses as the x variable that variable which is most highly correlated with the y variable. If the value of the most highly correlated variable is missing from the unit, the procedure uses the next most highly correlated variable. If that value is also missing, then the procedure continues in the same fashion until a reported value is found. Correlations are estimated from only the complete sample units.

This study uses two variations of the ratio procedure. These two variations arise because of the linear restriction imposed on the two variables -- $EXPFI + EXPFII \leq \text{Total Sows}$. The first variation simply imputes independently for EXPFI and/or EXPFII and then checks to see whether $EXPFI + EXPFII \leq \text{sows}$. When the total farrowings is greater than sows reported then the procedure proportionally adjusts any imputed values so that $EXPFI + EXPFII = \text{sows}$. The second variation uses the constructed variable $z = EXPFI + EXPFII$ as though it is a survey variable. If z is missing (either EXPFI or EXPFII is missing), the procedure: 1) finds an x variable by using correlations with z, 2) imputes a value for z, 3) makes $z = \text{sows}$ if $z > \text{sows}$, and 4) imputes for missing values of EXPFI and/or EXPFII so that $EXPFI + EXPFII = z$. If both EXPFI and EXPFII are missing, z is split into EXPFI and EXPFII proportionally by using relationships from the complete units in the data set.

2.2 The Array Procedure

The array procedure is not a procedure in general use but a procedure designed within ESS in 1971 [1] and proposed as a method of imputing for missing values on the multiple frame hog survey. Although not designed by the authors, the array procedure is included among the test procedures because its effects have never been assessed.

The array procedure uses a two-way table to impute for missing values. Two survey variables, total hogs and total sows, are chosen to define the table. If these two variables have c_1 and c_2 classes respectively, then the array procedure would form a table for EXPFI (as an illustration) of the form:

		Total Hogs			
		1	2	...	c_2
Total Sows	1				
	2				
	:				
	c_1				

Cell values must be initialized with an estimate for the farrowing ratio, $r = \text{EXPFI} / \text{the total number of sows}$. As the procedure processes the units in a sample, each unit is classified into a cell of the table by the reported values of the total hogs and the total sows for that unit. If EXPFI is reported, the value of r from the unit is added into a cell by using the weighted formula:

$$\frac{2 (\text{previous value for the cell}) + r}{3}$$

The purpose of this weighted formula is to prevent the imputation of extremely large values, i.e. outliers. If EXPFI is missing, the number of sows reported by the unit is multiplied by the ratio from the appropriate cell, and the result is imputed for the value of EXPFI. Obviously, the ordering of the data has some importance for the estimates from the array procedure. Although data from surveys by ESS - Statistics are often in a roughly geographic order, the data of this study were ordered randomly except that complete units were processed before complete units.

The array procedure is similar to the ratio procedure because the array procedure also used a type of ratio to impute values. However, the array procedure is a more complex method of obtaining the ratio and a more rigid process. For example, the array procedure uses the information from two other variables -- total hogs and total sows. However, these two variables are not allowed to have missing values. Another difference is that once the array procedure processes all complete units in a data set, then the procedure can also use incomplete units to change the ratio values in the cells as long as total hogs, total sows, and EXPFI are not missing. The ratio procedure, as used in this study can only use estimates of ratios and correlations from the complete units in a data set.

2.3 The ESTMAT Procedure

The ESTMAT procedure is an iterative solution to the problem of finding the maximum likelihood estimates for a multivariate data set in which some values are missing [4]. The ESTMAT procedure imputes by using multivariate regressions estimated from the reported values. As long as the same regression relationships apply to both reported and missing values, the ESTMAT procedure should be able to impute accurately even if the reported and missing values have different means.

The ESTMAT procedure represents an extension of the double sampling regression estimator to a multivariate setting. However, the ESTMAT procedure can take into account many different patterns of missing values in the data set. For example, once the data is collected for two variables, there are four possible patterns of missing data -- both variables are reported, only the first variable is reported, only the second variable is reported, or both variables are missing. With k variables there are 2^k possible patterns if one also counts as a pattern the set of complete units.

The estimation formulas which the ESTMAT procedure uses are complex and are not given in this paper. However, they can be found in the references [3] and [4]. Convergence of the iteration process used by ESTMAT is not assured in general, but in practical applications the convergence has usually taken less than ten iterations.

The two major assumptions of the ESTMAT procedures are: 1) values follow a multivariate normal distribution, and 2) the values are missing at random. The first assumption is necessary, of course, for the derivation of the maximum likelihood estimators used in the ESTMAT procedure. One example to show robustness to the normality assumption has been given [5], but no one has made a thorough study. The second assumption is unlikely to hold when the causes of the missing items are refusals, inaccessibles, editing, etc. The second assumption emphasizes the fact that the ESTMAT procedure seems more appropriate for survey situations in which the missing values are planned -- double sampling schemes, triple sample schemes, etc. [3]. However, if the procedure is robust to the randomness assumption, then applying multivariate regressions seems as reasonable as applying the ratio of a ratio procedure. The data set in this study does not obey either of the two assumptions for the ESTMAT procedure.

The ESTMAT procedure was *not* initially designed to impute individual values but to estimate directly the mean vector of the population. However, the procedure also estimates the variance-covariance matrix, and this estimate allows the computation of multivariate regression equations which can be used to impute individual values. These imputed values lack what Pregiborn [6] calls "commutativity" with the estimated mean vector. In other words, if one averages the reported and imputed values in a data set, this average does not equal the mean estimated directly by the ESTMAT procedure. Thus, the reader must be aware that the results of the ESTMAT procedure in this study are affected by an imputation process which may not be a part of other ESTMAT applications.

2.4 The Zero Spike Procedure

The zero spike procedure takes its name from the fact that zeros often dominate the response space of many surveys -- thus resulting in a "spike" of zeros when one draws a histogram of the distribution. The data set of this study has this characteristic. The first bar in each of the graphs of Figure 1.1 represents the zeros in the data set. For EXPFI 33 percent of the 201 original values are zeros, and for EXPFI 38 percent of the original values are zeros.

The zero spike procedure forms an indicator vector for each unit in the sample. For each variable, there is an element in the vector. The value of this element is "0" if the value of the variable is zero, "1" if the value is positive, and "2" is changed to a "0" or "1" using probabilities based on S -- that subset of the sample units which: 1) is complete, and 2) matches the indicator vector of unit A for those variables reported on A. For example, if there are two variables, then the complete units can form four groups -- (0,0), (0,1) (1,0) and (1,1). If a unit has the form (0,2) then the "2"

is changed to a "0" with probability $\frac{n(0,0)}{n(0,0) + n(0,1)}$ or changed to a "1" with

probability $\frac{n(0,1)}{n(0,0) + n(0,1)}$ where $n(i,j)$ is the number of complete units in

the (i,j) group; i, j=0,1. If a "2" is changed to a "0", the missing value becomes zero. If a "2" is changed to a "1", the missing value becomes a positive number of the form Rx , where x is the most highly correlated variable which also

has a "1" in the indicator vector of A. R is the ratio which relates x to y and is estimated from the units in S.

Pregiborn actually recommends the use of any, even subjective, information to estimate the probabilities for assignments of "0" and "1" and not just the use of units in the sample. Thus, his recommendations allow a Bayesian approach to the imputation through the estimation of the probabilities. Also, Pregiborn notes that there are many possible methods -- hot decks, regressions, averages, etc. -- to decide what positive value to impute for a missing value. This study uses a ratio method because the first three procedures described also use a ratio or regression method in some way. Thus, in the comparisons of estimates from the procedures, any differences for the zero spike procedure are not mainly a result of the method used to determine positive values but mainly a result of the "zero-positive" structure employed.

2.5 The Princomp Procedure

This procedure uses the first principal component when imputing for missing values. The first principal component is applied as a distance measure to select the complete unit which is most like a unit with a missing value. The reported value for this complete unit is then substituted for the corresponding missing value. The first principal component is a linear combination of all reported variables and has the maximum variance of all possible linear combinations of these variables. It is the line of closest fit in the sense that it minimizes the sum of squares of distances from data points to the line (note that a regression line minimizes the sum of squares in particular directions).

For this study the princomp procedure: 1) constructs four subsets of the data -- S1 contains the complete units, S2 contains those units with the variable EXPFI missing, S3 contains those units with the variable EXPFII missing, and S4 contains those units with both variables EXPFI and EXPFII missing; 2) computes the first principal component for S2 by using all 15 variables except EXPFI and then computes the value of the first principal component for each unit in S1 and S2; 3) for each unit in S2, finds the S1 unit which has a principal component value closest (minimum absolute deviation) to the unit in S2 and substitutes the corresponding values of EXPFI from the S1 unit into the missing values of EXPFI in the S2 unit; 4) repeats steps 2 and 3 to substitute reported values from S1 for missing values in S3 and S4 by using the principal component that corresponds to each subset.

The princomp procedure is essentially a hot deck procedure (a hot deck procedure is defined as a procedure which substitutes reported values for missing values) which substitutes by the minimization of a distance function rather than substituting randomly. There are many distance measures which could have been tested, but the authors felt that only one procedure of this type could be added to the experiment due to time and cost constraints and that the princomp procedure is a distribution-free method which has the potential for accurate imputation.

3. ANALYSIS

The goal of this analysis is to identify the "best" procedure of the six described procedures which impute for missing items. There are five criteria for selection of the "best" procedure: 1) the accuracy of estimated means,

2) the standard errors, 3) the accuracy of imputations on a unit level, 4) the effect on correlations between variables, and 5) costs.

3.1 Experimental Design

Three methods designate units which have missing items: 1) a random designation, 2) a 15 percent designation of incomplete units below the median and 85 percent above, and 3) an 85 percent designation of incomplete units below the median and 15 percent above. (The median of $z = \text{EXPFI} + \text{EXPFI}$ is used in these designations.) For each of these three methods, there are two rates to designate how many units have missing items -- 10 percent and 30 percent. The combined effect of the type of designation and the rate of designation results in six different situations in which means are estimated for the entire population.

Five data sets are simulated for each level of bias. Thus, a total of 30 data sets are generated from the original data set. Each data set consists of three replicates, and each procedure is run independently on each replicate to provide unbiased estimates of the standard errors. Within each data set the group of units with missing values contains 40 percent of the units with EXPFI missing, 40 percent with EXPFI missing, and 20 percent with both EXPFI and EXPFI missing.

The structure of the simulations corresponds to an analysis of variance model. This model is:

$$y_{ijklm} = \mu + \alpha_i + \beta_j + \gamma_k + T_\ell + (\alpha\beta)_{ij} + (\alpha T)_{i\ell} + (\beta T)_{j\ell} + \epsilon_{ijklm}$$

where:

α_i = the effect of the i^{th} designation method

β_j = the effect of the j^{th} rate of designating how many units have missing items

γ_k = the effect of the k^{th} replicate

T_ℓ = the effect of the ℓ^{th} missing item procedure

$(\alpha\beta)_{ij}$ = the interaction between designation method and rate

$(\alpha T)_{i\ell}$ = the interaction between designation method and missing item procedure

$(\beta T)_{j\ell}$ = the interaction between rate and missing item procedure

ϵ_{ijklm} = the error of the model associated with y_{ijklm}

y_{ijklm} = the value of the dependent variable associated with designation method i , rate j , replicate k , missing item procedure ℓ , and observation m

$i = 1, 2, 3$

$j = 1, 2$

$k = 1, 2, 3$

$\ell = 1, 2, \dots, 6$

$m =$ a number which varies with the definition of the dependent variable.

An example of a dependent variable is the difference for a sample unit between the imputed value and the corresponding original value. If an analysis of variance shows a significant difference due to an effect, then Duncan's multiple range test is used to identify which levels of the effect caused the differences. All tests are at a five percent level of significance.

3.2 Results

An analysis of variance shows significant differences among the six missing item procedures when the dependent variable is the average difference between the imputed values and the "true" values. Table 3.2.1 gives the results of Duncan's multiple comparison test and the patterns that are characteristic of each procedure. The ratio 1 and ratio 2 procedures are usually significantly different from the other procedures but not from each other. The princomp and zero spike procedures also tend to be different from the other procedures but are not significantly different from each other. The array procedure does not show consistent trends but tends to group with the princomp and zero spike procedures. The ESTMAT procedure tends to be by itself. Apparently the ESTMAT procedure is not robust to its normality and random error assumptions because the estimated means from this procedure are not very accurate under the random designation of missing values. All procedures tend to underestimate the mean -- even when values are randomly missing. This underestimation may not only be a result of biases inherent in the procedures but also a result of the skewness in the underlying data.

Table 3.2.1: Results of Duncan's multiple range test* when the dependent variable is the average difference between the imputed value and the corresponding original value.

Variable	Designation Method					
	Random		15% Below Median/85% Above		85% Below Median/15% Above	
	Average Difference	Procedure	Average Difference	Procedure	Average Difference	Procedure
EXPF1	-0.133	Ratio 2	3.781	ESTMAT	5.781	ESTMAT
	-1.281	Array	-5.719	Ratio 2	1.315	Ratio 2
	-2.359	Ratio 1	-5.922	Ratio 1	0.041	Ratio 1
	-5.285	ESTMAT	-10.715	Array	-4.104	Zero Spike
	-5.933	Zero Spike	-14.170	Princomp	-4.337	Princomp
	-6.756	Princomp	-15.870	Zero Spike	-5.759	Array
EXPF11	-0.852	Ratio 2	-7.567	Ratio 2	3.463	Ratio 2
	-1.500	Array	-8.711	Ratio 1	1.204	Ratio 1
	-2.104	Ratio 1	-14.378	Array	-2.641	Princomp
	-5.156	ESTMAT	-17.219	Princomp	-3.552	Array
	-5.815	Zero Spike	-18.330	Zero Spike	-4.285	Zero Spike
	-6.026	Princomp	-24.748	ESTMAT	-12.189	ESTMAT

* Any two means connected by the same bracket are not significantly different at $\alpha = 0.05$.

The interaction between the designation methods and the procedures is a significant effect. However, as Table 3.2.1 shows, this significance is a result of the fluctuation of the ESTMAT procedure in relation to the other procedures. The remaining tables in this paper give overall results across designation methods and rates. These overall results do not imply that the interactions are insignificant, but, as in Table 3.2.1, they are not important enough in this study to warrant the complexity of presenting the results in each cell. Table 3.2.2, for example, is much simpler and clearer than Table 3.2.1 and does not lose much information.

Table 3.2.2 gives overall results for Duncan's multiple comparison test in terms of average difference and relative bias. In this table the relative bias is the average difference in imputed and original values divided by the "true" mean of the sample. Across both variables the ratio 1 and ratio 2 procedures give the best results. It is disturbing that the ESTMAT procedure can give the best results for EXPFI and the worst for EXPFII. This result may be an effect of the imputation part of the ESTMAT procedure since direct estimates from ESTMAT showed a relative bias of -1.3 percent and -0.2 percent for EXPFI and EXPFII when estimating the mean for the entire population -- a result which seems more reasonable. Thus, *imputations* using the ESTMAT procedure appear to be unreliable.

Table 3.2.2: Overall results of Duncan's multiple comparison test*.

Variable	Procedure	Average Difference in Imputed Values and Original Values	Effect on Mean Estimates of Entire Population (Relative Bias)
EXPFI	ESTMAT]	1.426	+0.3%
	Ratio 2]	-1.512	-0.3%
	Ratio 1]	-2.747	-0.6%
	Array]	-5.918	-1.2%
	Princomp]	-8.421	-1.7%
	Zero Spike]	-8.636	-1.7%
EXPFII	Ratio 2]	-1.652	-0.4%
	Ratio 1]	-3.204	-0.8%
	Array]	-6.477	-1.3%
	Princomp]	-8.629	-1.8%
	Zero Spike]	-9.477	-2.0%
	ESTMAT]	-14.031	-2.9%

* Any two means connected by the same bracket are not significantly different at $\alpha = 0.05$.

To judge the accuracy of the imputations at a unit level, Table 3.2.3 gives the total of the absolute differences between the imputed values and "true" values. The optimum procedure should minimize this total. Table 3.2.3 confirms the superiority of the ratio 1 and ratio 2 procedures and explains the contradictory results in Table 3.2.2 between EXPFI and EXPFII for the ESTMAT procedure. The ESTMAT procedure gives the lowest difference for first quarter expected farrowings because of offsetting extremes in positive and negative directions. Thus, when absolute differences are calculated, the ESTMAT procedure gives the largest totals for both variables in Table 3.2.3.

Table 3.2.3: Total of the absolute differences between each imputed value and the corresponding "true" value.

Variable	Procedure	Absolute Difference
EXPFI	Ratio 2	6,683
	Ratio 1	6,711
	Array	9,648
	Zero Spike	9,909
	Princomp	10,129
	ESTMAT	14,643
EXPFII	Ratio 2	7,132
	Ratio 1	7,439
	Array	9,862
	Princomp	10,873
	Zero Spike	10,922
	ESTMAT	14,137

Table 3.2.4 gives the coefficient of variation of the estimated mean for the entire population. The coefficient of variation is the standard error (an unbiased estimate calculated using replicates) for a procedure divided by the "true" mean of the sample. The coefficients of variation in Table 3.2.4 are similar in size except that the ESTMAT procedure is larger for first quarter expected farrowings.

Table 3.2.4: Coefficients of variation for the estimated mean of the entire population.

Variable	Procedure	Coefficient of Variation
EXPMI	Zero Spike	0.062
	Princomp	0.063
	Array	0.065
	Ratio 1	0.065
	Ratio 2	0.070
	ESTMAT	0.100
EXPMII	Princomp	0.065
	Ratio 1	0.065
	Zero Spike	0.068
	Ratio 2	0.068
	ESTMAT	0.070
	Array	0.070

An overall measure of the quality of the procedures is the root mean square error. This measure is defined as:

$$\sqrt{MSE'} = \sqrt{(\text{Relative Bias})^2 + (\text{Coefficient of Variation})^2}$$

The $\sqrt{MSE'}$ is sensitive to the sample size since the sample size affects the magnitude of the coefficient of variation and sometimes the magnitude of the relative bias. Assuming, however, the relative bias is not affected by the sample size, Table 3.2.5 displays $\sqrt{MSE'}$ for several sample sizes by using the relative biases in Table 3.2.2 and the coefficients of variation in Table 3.2.4. Only for sample sizes larger than 1000 does the relative bias component dominate the root mean square error rather than the component due to the coefficient

of variation. Thus, for very large sample sizes, such as those often used in government surveys, the two ratio procedures give the best results. For smaller sample sizes, however, there is little difference in the procedures except that the ESTMAT procedure is substantially larger for EXPFI.

Table 3.2.5: Root mean square error relative to the "true" sample mean.

Variable	Procedure	Sample Size				
		50 (%)	100 (%)	1,000 (%)	10,000 (%)	∞ (%)
EXPFI	Ratio 1	13.0	9.2	3.0	1.0	0.6
	Ratio 2	14.0	10.0	3.1	1.0	0.3
	Array	13.1	9.3	3.2	1.5	1.2
	Princomp	12.7	9.1	3.3	1.9	1.7
	Zero Spike	12.5	8.9	3.3	1.9	1.7
	ESTMAT	20.1	14.2	4.5	1.4	0.3
EXPFI I	Ratio 1	13.0	9.2	3.0	1.2	0.8
	Ratio 2	13.6	9.6	3.1	1.0	0.4
	Array	14.1	10.0	3.4	1.6	1.3
	Princomp	13.2	9.4	3.4	2.0	1.8
	Zero Spike	13.8	9.8	3.6	2.2	2.0
	ESTMAT	14.3	10.3	4.3	3.1	2.9

When a data set contains imputed values, estimates of standard errors are often calculated by ignoring the imputation process and treating the imputed data set as though all the values are reported. This method may lead to biases in the estimates of standard errors. Table 3.2.6 gives the ratio of the variance calculated by using the conventional formula and the variance calculated by using replicates. The conventional formula treats the imputed values as though they are original, reported values. For example, in a simple random sample:

$$v_c(\bar{y}) = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n(n-1)}$$

where y_i is the value for observation i , $\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$, and n is the number of observations both reported and imputed in a simulated data set. An unbiased estimate of variance can be calculated by using replicates:

$$V_r(\bar{y}) = \frac{\sum_{i=1}^{n'} (y_i' - \bar{y}')^2}{n'(n' - 1)}$$

where y_i' is the mean of replicate i , $\bar{y}' = \frac{\sum_{i=1}^{n'} y_i'}{n'}$, and n' is the number of replicates. Table 3.2.6 shows there can be large biases in either direction for almost any of the procedures when using $V_c(\bar{y})$ as an estimate of the standard error. However, given the secondary importance of standard error estimates in the operational program as compared to total or mean estimates, the biases in the standard error estimates would probably not be very serious for the operational program of ESS.

Table 3.2.6: Ratio of estimated variances of estimated means -- variance estimate assuming imputed values are reported values divided by unbiased variance estimate using replication.

Variable	Designation Method	Imputation Procedure					
		Ratio 1	Ratio 2	Array	Zero Spike	Princomp	ESTMAT
EXPF1	Random	0.922	1.049	0.967	0.867	0.806	1.057
	15% Below Median/ 85% Above	1.084	1.172	0.970	0.902	0.746	1.316
	85% Below Median/ 15% Above	1.283	1.242	1.172	1.217	1.103	1.387
	Overall	1.096	1.154	1.036	0.995	0.885	1.253
EXPF11	Random	0.889	0.961	1.226	0.933	0.819	1.009
	15% Below Median/ 85% Above	0.869	0.933	0.980	1.063	0.838	1.136
	85% Below Median/ 15% Above	1.244	1.325	1.233	1.330	1.262	1.144
	Overall	1.000	1.073	1.146	1.109	0.973	1.096

Another important aspect of imputation is the effect on the correlation structure of the data set. Although correlations are not important for estimates of univariate statistics such as means and standard errors, correlations are important when the data set is used to explore and assess relationships among variables through regression analysis, principal components, or other multivariate techniques. Table 3.2.7 gives an example of the effects of the missing item procedures on the correlation structure. This table shows the correlations between sows and expected first quarter farrowings and between sows and expected second quarter farrowings. Most of the procedures tend to lower the correlations, but the ratio 1 and ratio 2 procedures tend to inflate the correlations.

Table 3.2.7: Correlations between SOWS and EXPFI and between SOWS and EXPFII for six missing item procedures.

Procedure	Variable					
	EXPFI			EXPFII		
	Random	15% Below Median/ 85% Above	85% Below Median/ 15% Above	Random	15% Below Median/ 85% Above	85% Below Median/ 15% Above
<i>(Actual)</i>	.82	.79	.84	.81	.72	.77
Ratio 1	.97	.94	.97	.89	.66	.93
Ratio 2	.88	.86	.94	.89	.72	.94
Array	.57	.73	.83	.54	.62	.76
Zero Spike	.80	.53	.74	.67	.42	.62
Princomp	.72	.60	.77	.68	.28	.74
ESTMAT	.79	.59	.81	.72	.33	.20

The cost of each procedure for imputing data is shown in Table 3.2.8. This is based on imputation for all 30 data sets for the two variables EXPFI and EXPFII. The system resource units (SRU's) -- a measure of computer usage -- required by each procedure are reasonably close except for the ESTMAT procedure. ESTMAT requires more SRU's than the other five procedures combined. This requirement is because of the complexity of the procedure. Thus, cost alone imposes a severe restriction on the use of the ESTMAT procedure. The other five imputation techniques are very similar in cost with the ratio 1 and ratio 2 procedures costing the least.

Table 3.2.8: Processing costs of six missing item procedures.

Procedure	SRU's ^{1/}	Cost ^{2/}
Ratio 1	937	\$ 145
Ratio 2	938	\$ 150
Array	1278	\$ 205
Zero Spike	1215	\$ 194
Princomp	1103	\$ 177
ESTMAT	9604	\$1537

^{1/} SRU: System resource unit

^{2/} Cost projected at 16¢ per SRU

4. SUMMARY

This study compares the effects of six procedures which impute for missing items -- two versions of the ratio procedure, the array procedure, the ESTMAT procedure, the zero spike procedure, and the princomp procedure. The comparison of these procedures is from an experiment in which a complete data set from a multiple frame hog survey by ESS-Statistics has values deleted to simulate an incomplete data set. Simulations are over a range of conditions which account for the method of designating missing values and the percentage of missing values. Comparisons of the procedures are made with respect to: 1) the accuracy of the estimated means, 2) the standard errors, 3) the accuracy of imputations on a unit level, 4) the effect on correlations between variables, and 5) costs.

The two versions of the ratio procedure perform the best for very large sample sizes (at least as large as 1000). For smaller sample sizes all of the procedures except the ESTMAT procedure have approximately the same mean square error. The main disadvantage of the ratio procedure is an inflation of the correlations between variables in the data set.

The ESTMAT procedure emerges as the least attractive procedure because it does not impute very accurately and it has an extremely high cost relative to the other procedures. This result only applies to the ESTMAT procedure as an *imputation* process and not as a missing data procedure in general. For example, the ESTMAT procedure is probably a suitable method for sample designs in which missing data is planned -- in other words, a survey design in which one plans to collect only partial information on some designated units.

Standard errors of the estimates from any missing item procedure should account for the fact that data is imputed. Replication is a method to obtain unbiased estimates of standard errors. Estimates of standard errors which treat the imputed values as though they are original, reported values may be biased, but the size of the bias is probably not very serious for the operational program of ESS.

Finally, the reader is cautioned that the results of the study are based on one data set in which the variables have skewed distributions dominated by zero values. These distributions are characteristic of much survey data collected by ESS but not all. Thus, generalizations of the results in this paper to other situations should be made with careful attention to the types of variables and their distributions and to the correlations among the variables.

5. RECOMMENDATIONS

Further research on the problem of missing items in the multiple frame hog survey should center on the operational procedure currently used by ESS. This procedure involves hand imputations for missing items by each field office. There are three aspects to this research: 1) documentation of the amount of hand imputing which is done, 2) measurement of how much this hand imputing affects the estimates, and 3) comparisons between the estimates from hand imputing and the other procedures described in this report, especially the ratio procedure. Because these comparisons would be under survey conditions, no "true" estimate will exist. Thus, it will be impossible to definitely state which procedure is the closest to "truth". However, relative comparisons of the estimates from the procedures can still provide insight into evaluating the problem of missing items and its solution.

BIBLIOGRAPHY

- [1] Beller, Norman D. and Bynum, Hugh E. "Multiple Frame Hog Survey; Nebraska's Hot Deck Edit Procedure, Version 2." Documentation of computer program. 1971.
- [2] Ford, Barry L. "Missing Data Procedures: A Comparative Study", Proceedings of the Social Statistics Section, pages 324-329. American Statistical Association. 1976.
- [3] Ford, Barry L.; Hocking R. R.; and Coleman, Ann. "Reducing Respondent Burden on an Agricultural Survey", Proceedings of the Section on Survey Research Methods, pages 341-345. American Statistical Association. 1978.
- [4] Hartley, H. O. and Hocking, R. R. "The Analysis of Incomplete Data", Biometrics. Volume 27, pages 783-823. 1971.
- [5] 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, pages 121-133. 1974.
- [6] Pregiborn, Daryl. "Incomplete Survey Data: Estimation and Imputation", Methodology Journal of Household Survey Division. 1976.
- [7] Pregiborn, Daryl. "Discussion of Papers by Huddleston and Hocking and Patrick", Proceedings of the Section on Survey Research Methods, pages 492-493. American Statistical Association. 1978.