# A STUDY OF DONOR POOLS AND IMPUTATION METHODS FOR MISSING EMPLOYMENT DATA

Kenneth W. Robertson, Albert Tou, Larry Huff
Keneth W. Robertson, Suite 4985 Bureau of Labor Statistics, 2 Massachusetts Ave. NE, Washington D.C., 20212

## Outline

## I. General Information

The Bureau of Labor Statistics' Occupational Employment Statistics (OES) survey is a periodic survey of nonfarm establishments that collects occupational employment data on workers by industry. The survey uses a weighting cell adjustment procedure to adjust for unit nonresponse. Previous research has shown that this procedure works reasonably well with these data. However, a weakness of this procedure, as employed in this survey, is that it does not adjust for unit nonresponse in three-digit industries that have no responding units. The adjustment cell for the current procedure is, at most, an entire three-digit SIC within a sampled area.

The OES sample is originally allocated at the Substate-Area, three-digit SIC, Employment-Size level. Estimates are produced at the Substate-Area, three-digit SIC level. Estimates at higher levels of aggregation are then summed from these lower level estimates. Estimates are generally published by each state at the Substate-Area, two-digit SIC level. Estimates at lower levels, e.g. three-digit SIC levels, are generally used internally as an input in the production of occupational employment projections. For most, if not all states, a primary reason to conduct the survey is to produce occupational employment projections. Therefore, it is important that estimates at both the two-digit and three-digit SIC levels be as accurate as possible.

The impetus for this research was observations from users about the data processing system not handling situations where there were no data in the three-digit industry / area cell (empty industry / area cell) to use for the weighting cell adjustment. The result is that the two-digit SIC / area estimates were incomplete in that they do not include the empty three-digit industry / area cell and then had to be described in the published materials.

## II. Data

As a first step in conducting this research we obtained data from three States, for three consecutive years. Since the OES is a survey which takes three years to cycle through all the sampled industries, this gave us a sample of all available industries in three States. The data contained three stratification variables, each with three levels. These variables, and their associated levels, are as follows:

1. Location

| Substate-Area | Met | State-Wide |

2. Industry

| three-digit (3D) SIC | Near-SIC | two-digit SIC |

3. Employment Size

| Sampled-Size | Near-Size | All-Size |

Met was derived from auxiliary data and indicates whether or not the establishment is located within a metropolitan area. Near-SIC includes reporting units in both the sampled three-digit SIC and the three-digit SICs on either side numerically, as long as they remain within the two-digit SIC. No adjustment is made for non-existing three-digit SICs. Near-Size is similar, except that this would include the size classes to either side of the sampled size class.

Three typical occupations were selected for each two-digit SIC.

Response for the OES survey is generally above 75 percent, and sometimes approaches 100 percent in some States and three-digit SICs. The actual nonresponse within each sampling cell (Substate-Area, three-digit SIC, Sampled-Size) was calculated for our data. All nonrespondents were then removed from the dataset. Nonresponse was simulated by using the actual nonresponse rate in each stratum to randomly pick units from the respondents to act as nonrespondents. This was repeated 25 times, producing 25 data sets, each with its own set of respondents and randomly selected units to act as nonrespondents.

Note that these data sets are not independent, as each unit is represented in each dataset, the only variation among the 25 data sets is whether or not a given unit is posing as a respondent or as a nonrespondent. Using multiple simulations in this manner provides the opportunity to observe what happens to the estimates as different sets of units are deleted to simulate nonrespondents and are represented by the remaining respondents.

## III. Donor Pool Research

Our goal for this project is to determine whether there is an imputation method which is at least as good as the current weighting cell method, without its shortcomings, as it is applied in this survey. First, though, a determination of where to get the data with which to impute is needed. It is intuitive that the data should be taken from the sampling cell. Obviously, the units within the sampling cell are homogeneous with respect to the stratification variables. However, as was previously mentioned, there are instances where no data are available in a particular

cell to use for the weighting cell adjustment. In this event, some objective criteria are needed to determine which donor pools to use. A donor pool is a set of responding units within a defined set of strata. The first step in the research was an examination of this problem.

Multiple donor pools were explored by allowing the donor pool to expand in increments to all areas and multiple three-digit SICs. Later, alternative nonresponse adjustment procedures such as mean imputation, hot deck (nearest neighbor), and hot deck (random selection in a cell) are used to adjust for nonresponse utilizing the expanded donor pools in the event that the original pool is empty. The first non-empty pool encountered will be used. The accuracy of each of these procedures in producing aggregate level estimates is determined. Of particular interest is the accuracy of aggregate level estimates that include three-digit industries having no responding units.

The purpose in examining donor pools is to determine an appropriate expansion for those instances where imputation is not possible (because no data are available). What is needed is a method of expanding the donor pool from the level at which units were allocated to the maximum allowable level (State-Wide, two-digit SIC, All-Size). The expansion should occur in such a way as to arrive at the 'best' pool of donors available. The expansion will be utilized only when there are no data available in the sampled cell. 'Best' is defined as being that pool which, in general, results in estimates with smaller errors than the pool that follows it in the expansion.

To determine the order of the donor pool expansion, some imputation method was needed to generate estimates and errors for each of the possible donor pools. Since mean imputation is both easy to implement and should give results reasonably close (due to units changing SICs, weighting cell and mean imputation produce different estimates) to the current weighting cell method, it was determined to be appropriate for this purpose. Specifically, for each nonrespondent values for occupational employment were imputed using the mean occupational employment derived from the respondents within each of the 27 possible donor pools.

Once the imputation was completed, we had 25 datasets, each containing respondents and nonrespondents. Each nonrespondent contained 27 imputed values for each of the three typical occupations. Using the reported data from the respondents, and the imputed values from the nonrespondents, estimates were produced at the three-digit SIC / Area level. Therefore, for each three-digit SIC / Area in the survey, there were 27 estimates across 25 iterations for each of the three typical occupations. Given that there were 829 three-digit SIC / Areas, there were a lot of data to evaluate. At this point an appropriate expansion of cells was all that we were seeking. Therefore all cases which could not be imputed, and thus could not be used for comparisons of different donor pools, were eliminated. This eliminated 2.7 percent of the nonrespondent records, and resulted in removing 1.5 percent of the SIC / Area estimates.

In order to determine an expansion from these data some method of summarizing the data was needed. The data were summarized in two ways. First, the absolute error of the estimates was computed, and summarized as follows:

$$t_{c,O} = \sum_{Sz} \sum_{r+nr} w_{c,u} * rd_{c,u,O}$$

$$e_{c,O,P} = \sum_{Sz} \sum_{r} w_{c,u} * rd_{c,u,O} + \sum_{Sz} \sum_{nr} w_{c,u} * id_{c,u,O,P}$$

$$err_{c,O,P} = \left| t_{c,O} - e_{c,O,P} \right| \qquad rank_{c,O,P} = rank(err_{c,O,P}) \qquad \text{Where}$$

$$\{rank = 1,2,...,27, \text{ ties allowed}\}$$

$$\overline{rank}_{c,P} = \frac{1}{3} \sum_{O=1}^{3} rank_{c,O,P} \qquad \overline{\overline{rank}}_{P} = \frac{1}{n_c} \sum_{c} \overline{rank}_{c,P}$$

$$total\_err_P = \sum_{c} \sum_{O} err_{c,O,P}$$

Where
  c = cell : Year, State, three-digit SIC, SubState Area
  P = donor Pool {1,2,...,27}
  Sz = Employment Size Class , O = Occupation {1,2,3}
  r = respondents , nr = nonrespondents
  u = establishment , $n_c$ = number of cells
  $w_{c,u}$ = the weight for establishment u in cell c
    This is the inverse of the probability of selection into the sampled cell
  $rd_{c,u,O}$ = reported occupational employment data for occupation O from establishment u within cell c
  $id_{c,u,O,P}$ = imputed data for occupation O from establishment u within cell c, data imputed using Pool P
  $t_{c,O}$ = true employment value within cell c for occupation O
  $e_{c,O,P}$ = estimated total employment for occupation O within cell c using reported data for respondents and data imputed from Pool P for nonrespondents
  $err_{c,O,P}$ = the absolute difference of the true value and the estimated value
  $rank_{c,O,P}$ = the rank of $err_{c,O,P}$ within each cell across the donor pools.

Ties received the same rank, thereby reducing the maximum rank for that record. These ranks were averaged across the three occupations, and then averaged across all areas, SICs, and States. This provided a way to rank the donor pools at a global level. The absolute errors were also summed across all areas, SICs, and States. This produced a method of measuring the errors associated with each donor pool at a global level. The two sets of summary statistics, ranks and summed absolute errors, provided the information needed to determine an appropriate expansion scheme. To verify that the results were consistent, we also calculated errors at the unit level. As before the average of the ranks was computed, as well as the sum of the absolute errors.

Aggregate summaries were used for this phase of the research because the expansion of donor pools must be a general case expansion. While there are probably specific cases where other expansions might be more beneficial, any method chosen for the OES survey will have to work well in all cases, and as such can not be too specific in nature. The summaries chosen here should yield an expansion which works well in most cases.

## IV. Results from III

The results were quite interesting. First, it was very obvious that the most important donor pool stratification variable is Employment Size. The errors and ranks were both clearly grouped by the three values of this variable in the following order; Sampled-Size, Near-Size, and All-Size. There also appeared to be groupings by industry, also in order, however, this was not quite as clear as the size class groupings. The final variable, location, did not seem to make a large difference after accounting for Size and Industry.

In order to verify these results, we used the Cox-Stuart nonparametric test for Trend. We structured this test as follows. Our intuition tells us that within the Size stratum, the levels should be ordered from Sampled-Size to Near-Size to All-Size. Pairs where this order held were assigned a minus sign, those pairs where it failed were assigned a plus sign. This resulted in n pairs, where the number of plus signs is our test statistic k. This statistic is compared to a binomial table to determine its associated p-value. If the trend holds there will be very few plus signs, and a correspondingly small p-value.

This test was also done for the Industry variable, and the Location variable. Once we had determined the variable which was most important overall, we produced p-values for tests within each of the three subgroups formed by the most important variable. For example, if Size were most important, then we would test the values of Location and Industry which fell within the grouping Sampled-Size, then those values of Location and Industry which fell within the grouping Near-Size, and finally those values of Location and Industry which fell within the grouping All-Size.

Since we produced these rank orders in four ways, we conducted these tests four times. The results for the first test are given in the table below. The remaining three tests provided similar results in most cases. The p-values indicate the significance level of the test. We reject the null hypothesis of no trend if the p-value is less than 10 percent. Consequently the alternative hypothesis, that there is a trend in the direction chosen, is accepted.

| | Location | Industry | Size |
|---|---|---|---|
| All | p=0.6367 | p=0.1937 | p=0.0001 |
| Size=Sampled-Size | p=0.6875 | p=0.0625 | |
| Size=Near-Size | p=0.500 | p=0.0625 | |
| Size=All-Size | p=0.500 | p=0.0625 | |

As the Cox-Stuart tests indicate, there is a significant trend in the indicated direction (Sampled-Size, Near-Size, All-Size) for the Size variable at the p=0.0001 level. Additionally, within these groupings by size, the Industry variable relatively often shows a trend in the selected direction (three-digit SIC, Near-SIC, two-digit SIC) at the p=0.0625 level. There were three instances where reversing the direction would have yielded p-values less than 10 percent. The test for a trend in the Location variable was never significant in the direction we chose (Substate-Area, Met, State-Wide). However, there were four cases where reversing the direction would have yielded p-values less than 10 percent. These statistics support our findings.

Based on these results, the following order was chosen for the expansion of donor pools. The Met value of the location variable was eliminated, since it seemed to have little effect.

**Donor Pool Expansion Order**

| | Stratum | | |
|---|---|---|---|
| | Location | Industry | Employment Size |
| 1 | Substate-Area | three-digit SIC | Sampled-Size |
| 2 | State-Wide | three-digit SIC | Sampled-Size |
| 3 | Substate-Area | Near-SIC | Sampled-Size |
| 4 | State-Wide | Near-SIC | Sampled-Size |
| 5 | Substate-Area | two-digit SIC | Sampled-Size |
| 6 | State-Wide | two-digit SIC | Sampled-Size |
| 7 | Substate-Area | three-digit SIC | Near-Size |
| 8 | State-Wide | three-digit SIC | Near-Size |
| 9 | Substate-Area | Near-SIC | Near-Size |
| 10 | State-Wide | Near-SIC | Near-Size |
| 11 | Substate-Area | two-digit SIC | Near-Size |
| 12 | State-Wide | two-digit SIC | Near-Size |
| 13 | Substate-Area | three-digit SIC | All-Size |
| 14 | State-Wide | three-digit SIC | All-Size |
| 15 | Substate-Area | Near-SIC | All-Size |
| 16 | State-Wide | Near-SIC | All-Size |
| 17 | Substate-Area | two-digit SIC | All-Size |
| 18 | State-Wide | two-digit SIC | All-Size |

## V. Imputation Research

In the imputation research, three imputation methods were used, Hot-Deck (Nearest Neighbor), Hot-Deck (Random Selection within a cell), and Mean of Cell.

The Hot-Deck Nearest Neighbor was implemented as follows. When the sample is drawn, an employment value is taken from the frame. This employment value, denoted as the Original Benchmark Employment (OBME), is placed on the file and is used to place the unit within the appropriate size strata. This value exists on all establishment records, whether they have responded or not. The respondent within the cell that had an OBME value closest to the OBME value of the nonrespondent was located. This respondent's employment values were then used as the imputed employment values for the nonrespondent.

The Hot-Deck random selection within a cell method was implemented in the following way. A uniform random number was assigned to all units. The respondent within the cell that had a random number value closest to the random number value of the nonrespondent was chosen. This respondent's employment values were then used as the imputed employment values for the nonrespondent.

The mean of cell imputation method was implemented in the usual way. All respondents within the cell were used to find a mean employment value for each occupation. These mean occupational values were then used as imputed occupational values for each nonrespondent.

Estimates were also computed using the current method, which is Weighting Cell adjustment. This method applies a weight adjustment to a reporting establishments occupational data

which is equal to the weighted OBME of all establishments in the cell divided by the weighted OBME of the responding establishments.

These methods are referred to as NN (Nearest Neighbor within Cell, based on the frame employment value), RS (Random Selection within Cell), MI (Mean of Cell Imputation), and WC (Weighting Cell) in the tables presented later in the paper. All imputations were done at the sampling cell level, except when there were no data available at that level. When no data were available, the donor pool expansion order listed previously was followed until data became available. The weighting cell method was an exception to this. Since there was a need to compare imputation methods with the current method, the weighting cell method was computed as it is currently implemented. When no data are available at the sampled level, the current method uses an expansion alternating about the current size class until either an acceptable NonResponse Adjustment Factor (NRAF) is calculated, or all size classes have been included. There is no utilization of other SICs or Areas. If all size classes have been used and no data have been found then the current method fails, and we have an empty industry / area cell in the two-digit estimates when they are summed from the lower estimating level. The incidence of these empty industry / area cells was determined to be 9.1 percent in the original data. By allowing imputation across substate areas this empty cell rate drops to 0.6 percent. A further expansion to the State-wide / two-digit SIC level eliminates the problem entirely. Therefore, by incorporating an imputation method which allows data to be imputed from a wider pool than that currently employed in this survey, this problem could be eliminated.

The data used in the donor pool research were also used to test imputation methods. Twenty-five data sets were produced with nonrespondents randomly chosen from the set of available units. Again, these data sets were not independent, since each data set contained identical units. Only the assignment of respondent or nonrespondent changed across the data sets.

These data sets were used to impute for the nonrespondents using the four methods listed, for each of the three selected occupations. Estimates of total occupational employment and variances of mean occupational employment were computed using the imputed values for the nonrespondents and the reported values for the respondents. These estimates and variances were then compared against the estimates and variances generated using only reported data from each unit. This resulted in errors for each of the 829 State / SIC / SubState Area cells, for each of the three occupations, across each of the four imputation methods.

The estimates and errors were calculated as follows:

$$t_{c,O} = \sum_{Sz} \sum_{r+nr} w_{c,u} * rd_{c,u,O}$$

$$e_{c,O,I} = \sum_{Sz} \sum_{r} w_{c,u} * rd_{c,u,O} + \sum_{Sz} \sum_{nr} w_{c,u} * id_{c,u,O,I}$$

$$err_{c,O,I} = t_{c,O} - e_{c,O,I} \qquad \overline{m}_{O,I} = \sum_{c \in A} \frac{err_{c,O,I}}{n_A}$$

$\tilde{m}_{O,I}$ = that $err_{c,O,I}$ which satisfies :

$$P(err_{O,I} < err_{c,O,I}) \leq \frac{1}{2} \quad \text{and} \quad P(err_{O,I} \leq err_{c,O,I}) \geq \frac{1}{2}$$

$$v_{O,I} = \sum_{c \in A} \frac{(err_{c,O,I} - \overline{m}_{O,I})^2}{n_A - 1}$$

$$rank(|x|) = rank(|x|, \text{ across I})$$

$$(x = \overline{m}_{O,I}, \tilde{m}_{O,I}, v_{O,I}) \quad \left(\text{values} = \{1,2,3,4\}\right)$$

$$\overline{\overline{m}}_{SIC2D,I} = \sum_{FIPS} \sum_{O} \frac{1}{9}(rank(\overline{m}_{O,I})_{SIC2D,FIPS})$$

$$\overline{\tilde{m}}_{SIC2D,I} = \sum_{FIPS} \sum_{O} \frac{1}{9}(rank(\tilde{m}_{O,I})_{SIC2D,FIPS})$$

$$\overline{v}_{SIC2D,I} = \sum_{FIPS} \sum_{O} \frac{1}{9}(rank(v_{O,I})_{SIC2D,FIPS})$$

Where

$err_{c,O,I}$ = the estimated employment value for cell c, Occupation O, & Imputation method I minus the true employment value for cell c & Occupation O

$err_{O,I}$ = the distribution of errors, $err_{c,O,I}$

c = cell = Year, State, three-digit SIC, SubState Area

A = all c (cells) within the Statewide two-digit SIC

$n_A$ = the number of cells in set A

I = NonResponse Adjustment Method {NN, RS, MI, WC}

O = Occupation {1,2,3} , Sz = Employment Size Class

r = respondents , nr = nonrespondents

$t_{c,O}$ = true employment value within cell for occupation O

$e_{c,O,I}$ = estimated employment value within cell for occupation O

$\overline{m}_{O,I}$ = the mean error for occupation O when imputation method I is used, i.e., the error averaged across A

$\tilde{m}_{O,I}$ = the median error for occupation O when imputation method I is used, i.e., that error which falls in the middle when the errors are ordered from smallest to largest within A

$v_{O,I}$ = the variance of the imputation error for occupation O when imputation method I is used within A

rank(x) = the above three quantities were then ranked across I, so that each imputation method was given a value between 1 and 4, the imputation method resulting in a value closest to 0 being assigned a 1. The imputation method resulting in the value farthest from 0 was assigned a 4.

The final three statistics, $\overline{\overline{m}}_{SIC2D,I}$, $\overline{\tilde{m}}_{SIC2D,I}$, & $\overline{v}_{SIC2D,I}$, are averages of the indicated **ranks** across State and Occupation.

Errors ($err_{c,O,I}$) for both the estimate (of total employment-shown above) and variance (of mean employment) were computed at the three-digit SIC, SubState Area level. The mean, median, and variance of these errors were computed. The absolute values of these statistics were then ranked by imputation method to determine the best imputation method for that cell. These ranks were then averaged across States and Occupations to give us an

617

average mean, average median, and average variance rank for the errors from each two-digit SIC / Imputation method. Ties received the same rank, thereby reducing the maximum rank for that two-digit SIC. The results produced information about the bias of the estimates under each of these imputation methods. That is, the method with the smallest average rank generally has the least bias associated with that imputation method.

There was also interest in the absolute size of the errors, therefore, similar statistics for the absolute errors were computed as above with the absolute value of $err_{c,O,I}$. These results will tell us which imputation methods has the smallest errors, regardless of any bias which may be present.

Placing one or more values at the mean adds nothing to the variance, therefore the variance has a downward bias when utilizing the Mean of Cell imputation method. In previous research (West, Kratzke, and Robertson, 1994), an adjustment was derived to adjust for this downward bias, based on the response level. As shown in the previous research we can define the population variance as

$$\text{follows: } V_t = \frac{\sum\limits_{j \in r}\left(E_j - \overline{E}\right)^2 + \sum\limits_{k \in nr}\left(E_k - \overline{E}\right)^2}{n_r + n_{nr}}$$

Assuming that the nonrespondents are missing at random, we can consider the effect that the mean imputation method has on this variance. When mean imputation is used, the second summation becomes zero, since $E_k = \overline{E}$, giving us the following

$$V_i = \frac{\sum\limits_{i \in r}\left(E_i - \overline{E}\right)^2}{n_r + n_{nr}} = \frac{\left(n_r - 1\right)(S)^2}{n_r + n_{nr}}$$

$$\text{where } (S)^2 = \frac{\sum\limits_{j \in r}(E_j - \overline{E})^2}{n_r - 1}$$

where
   $V_t$ = the "true" variance, and
   $V_i$ = the variance when mean imputation is used
   $n_r$ = number of respondents, and
   $n_{nr}$ = number of nonrespondents.

Since $S^2$ is an unbiased estimator of $V_t$, $E(V_i) = \dfrac{n_r - 1}{n_r + n_{nr}} V_t$.

Therefore, we can remove the bias from $V_i$ by making the following adjustment: $V_t \approx \dfrac{n_r + n_{nr}}{n_r - 1} V_i$.

A variation of this adjustment has been applied to the variance estimates calculated for the Mean of Cell imputation method. The adjustment used is given in the following equation.

$$V_{Adj} = V_S * \frac{WE_R + WE_{NR}}{WE_R - 1}$$

Where
   $V_{Adj}$  is the adjusted variance estimate
   $V_S$ is the standard variance estimate
   $WE_R$  is the weighted frame employment of the respondents, and
   $WE_{NR}$ is the weighted frame employment of the nonrespondents

This adjustment produces an unbiased variance estimator to use in evaluating the Mean of Cell imputation method.

## VI. Results from V

Our first finding from this research is that we were able to impute for all nonrespondents after the fifth expansion of the donor pool to the State-wide / two-digit SIC / sampled size donor pool. As indicated earlier, if there were no respondents in the sampled pool, the donor pool underwent 17 expansions until we found data with which to impute. Our overall statistics show us that without any expansion we were able to have data for 93.5 percent of the units. This figure includes both respondents and units which had been imputed to this point. The first expansion gave us data for 99.4 percent of the units. After the third donor pool, we only had 0.2 percent of the units left to impute.

Errors for both the employment estimate and the variance were computed at the State / 3-digit SIC / Substate Area level for three occupations using four imputation methods. The mean, median and variance of these errors were computed at the State / 2-digit SIC level.

As shown in the previous section, the absolute values of these six statistics were then ranked, providing a rank (1-4) for each State / 2-Digit SIC / Occupation / Method. These ranks were then averaged across the states and over the three occupations, giving an average rank for each 2-Digit SIC / Method. Table I summarizes the results. Using ranks it is clear that the Nearest Neighbor imputation method is giving us the best results. It is also clear that the current method, Weighting Cell, is giving us the worst results. In many cases the differences in estimation errors were minute. For most SICs, the difference between the best estimate and the worst estimate is less than 1, which means they would be rounded to no difference. Even though the Nearest Neighbor came in first place in almost every category, any of these imputation methods would provide a good estimate. However, there were larger differences in the variance estimates. In this arena, the Nearest Neighbor shows a clear edge over the Weighting Cell.

The previous rankings provide information on the bias of the errors. Also of interest was which method produced the smallest errors, without regard to bias. To determine this the ranks of the absolute errors were examined. Table II shows that the Nearest Neighbor method again comes in first place, although the Weighting Cell does not do as bad as before. There is less difference among the ranks when the absolute errors are examined.

Our results show that the Nearest Neighbor is the least biased of the imputation method for both the estimate and variance, and that it produces the best absolute errors for the variance.

An examination of the distributions of estimation errors made it clear that there is generally little difference between the best and worst method. Therefore, any of the methods used in this research would be acceptable for estimation of these data. Also clear, however, is the fact that the distributions of variance errors were not alike. The distribution of variances for the Weighting Cell method shows more dispersion in almost every case than the Nearest Neighbor method. Based on these distributions we expect the Nearest Neighbor imputation method to produce a variance error which will generally be closer to zero than the Weighting Cell adjustment variance error.

## VII. Conclusions

Originally, we stated that we were particularly interested in analyzing those estimates which included empty Area / SIC cells. However, we also noted earlier that imputation would eliminate the problem. Because of this we did not measure the gains in precision due to imputation. It is sufficient to say that using the current methodology results in occasional empty Area / SIC cells, while the imputation methods examined in this paper eliminate this problem and allow two-digit estimates to be produced without referring to empty cells.

Our primary concern is to suggest a nonresponse adjustment method which will produce good estimates, good variances, and will also solve the occasional empty Area / SIC problem. Based on our results here, it seems that the Nearest Neighbor method of imputation will meet all of these requirements.

## VIII. Future Research

In the future we plan to continue searching for ways to improve the OES survey. One area of interest is the sample allocation methodology. The procedure currently in use was adapted from a procedure designed prior to widespread computer availability. It is felt that this process could be made more efficient by redesigning it to take advantage of technological advances. A second area of interest is to test alternative variance estimators, and see if any of these are more efficient than the current method.

## References
1. Sandra A. West, Diem-Tran Kratzke, and Kenneth W. Robertson (1994), "Variance Estimators For Variables That Have Both Observed and Imputed Values"
*ASA Proceedings of the Section in Survey Research Methods.*

**Table I  Bias Measurement**

How many times was the method the **best**?
(The average rank was the smallest)

|  | MI | NN | RS | WC |
|---|---|---|---|---|
| Mean of Estimate | 2 | 4 | 1 | 2 |
| Mean of Variance | 0 | 6 | 5 | 0 |
| Median of Estimate | 0 | 9 | 0 | 0 |
| Median of Variance | 0 | 9 | 0 | 0 |
| Variance of Estimate | 3 | 1 | 0 | 7 |
| Variance of Variance | 1 | 6 | 2 | 0 |
| Sum | 6 | 35 | 8 | 9 |

How many times was the method the **worst**?
(The average rank was the largest)

|  | MI | NN | RS | WC |
|---|---|---|---|---|
| Mean of Estimate | 2 | 3 | 3 | 2 |
| Mean of Variance | 1 | 0 | 0 | 8 |
| Median of Estimate | 0 | 0 | 0 | 9 |
| Median of Variance | 4 | 0 | 0 | 6 |
| Variance of Estimate | 0 | 3 | 6 | 1 |
| Variance of Variance | 0 | 0 | 0 | 9 |
| Sum | 7 | 6 | 9 | 35 |

**Table II  Absolute Error Measurement**

How many times was the method the **best**?
(The average rank was the smallest)

|  | MI | NN | RS | WC |
|---|---|---|---|---|
| Mean of Estimate | 1 | 1 | 0 | 7 |
| Mean of Variance | 0 | 7 | 2 | 0 |
| Median of Estimate | 2 | 6 | 0 | 1 |
| Median of Variance | 0 | 7 | 2 | 0 |
| Variance of Estimate | 4 | 0 | 0 | 6 |
| Variance of Variance | 1 | 5 | 3 | 0 |
| Sum | 8 | 26 | 7 | 14 |

How many times was the method the **worst**?
(The average rank was the largest)

|  | MI | NN | RS | WC |
|---|---|---|---|---|
| Mean of Estimate | 0 | 4 | 5 | 1 |
| Mean of Variance | 1 | 0 | 0 | 8 |
| Median of Estimate | 1 | 1 | 2 | 5 |
| Median of Variance | 5 | 0 | 0 | 4 |
| Variance of Estimate | 0 | 3 | 5 | 1 |
| Variance of Variance | 0 | 0 | 0 | 9 |
| Sum | 7 | 8 | 12 | 28 |