

I. Introduction

The Monthly Retail Trade Survey is a panel survey that provides retail sales information from a probability sample. The list sample for any given data month consists of certainty units from a fixed panel and sampling units from a rotating panel. The certainty units report the sales monthly; while the sampling units from the rotating panel report both current and previous month sales every three months.

The missing item of the current month retail sales is imputed by multiplying the nonresponding unit's previous month sales (reported or imputed) by a measure of trend (the so called ratio of identicals) computed from those responding units whose size and kind of business characteristics are similar to the nonresponding units. The trend is calculated based on the weighted sum of the current month sales and the weighted sum of the previous month sales for each adjustment cell. The sample is partitioned into imputation cells defined by kind of business (3-or-4 digit Standard Industrial Classification (SIC) code is used), firm size (Group I and Group II; Group I has 3 different firm sizes (1-3, 4-10, 11+ establishments), Group II is the firm with 11+ establishments and certainty alpha) and size of sales (defined arbitrarily). Missing previous month sales for a sampling unit in a rotating panel is imputed in the similar fashion using historical data.

In this paper, we compare several ratio and regression adjustment procedures and a variety of imputation cell formations under a Monte Carlo study. We treated data reported in the Monthly Retail Trade Survey from 9 SIC's as our complete data set, and assumed that the data were missing at random. The missing items were imputed by different imputation procedures. The bias and mean square error (MSE) of the estimated totals for the given data set are derived in the following section. The conclusion in the study will be summarized in section III.

II. Monte Carlo Study

In Huang (1984), a Monte Carlo study was carried out to evaluate different imputation procedures based on a given set of complete data (reported list sample from SIC 562 in the December 1982 Retail Trade Survey). Five sets of incomplete data were generated from the complete data. For each of the five sets, data were randomly suppressed from each imputation cell of complete data according to its current imputation rate. The reader is cautioned that since only five incomplete data sets were used, the results of the comparisons may not give an accurate picture. In the following, the bias and MSE of the estimated total for a given complete data set were derived under the assumption that the missing data are a random sample of the complete data set.

A sample of size  $n$  is assumed to be drawn from a population of size  $N$ . Only one imputation cell is assumed. The sampling unit  $i$  has inclusion probability  $\pi_i$ . In a sample of size  $n$ , there are  $n_r$  units reported, and  $n'_r$  units not reported. In the following we treat these reported  $n_r$  units as our complete data set. Assuming the nonresponse mechanism is ignorable, i.e., the data are missing at random, the incomplete data sets are generated in which  $n'_r$  units

are suppressed randomly from the complete data set of  $n_r$  units, and different ratio type imputation procedures are used to impute  $n'_r$  missing units of  $y$  values (current month sales) using the auxiliary variable  $x$  (previous month sales), which is available for all  $n_r$  units.

Let  $Y$  be the estimated total using the complete data set of  $n_r$  units

$$Y = \sum_{i=1}^{n_r} y_i / \pi_i.$$

Let  $\hat{Y}$  be the estimated total using incomplete data set of  $n_r$  units, of which  $n'_r$  units are reported, and  $n''_r$  units are imputed, i.e.,

$$\hat{Y} = \sum_{i=1}^{n'_r} y_i / \pi_i + \sum_{i=1}^{n''_r} \hat{y}_i / \pi_i$$

where

$$\hat{y}_i = \hat{R}_{n'_r} x_i, \quad i=1, \dots, n''_r,$$

$$n_r = n'_r + n''_r,$$

$\hat{R}_{n'_r}$  is one of the four ratio type estimators in (2.1)-(2.4) using  $n'_r$  units.

$$\hat{R}^{(1)} = (\sum_i y_i / \pi_i) / (\sum_i x_i / \pi_i) \quad (2.1)$$

$$\hat{R}^{(2)} = \sum_i y_i / \sum_i x_i \quad (2.2)$$

$$\hat{R}^{(3)} = (\sum_i \frac{1}{\pi_i} \frac{y_i}{x_i}) / (\sum_i \frac{1}{\pi_i}) \quad (2.3)$$

$$\hat{R}^{(4)} = \frac{1}{n'_r} \sum_{i=1}^{n'_r} (y_i / x_i) \quad (2.4)$$

$\hat{R}^{(1)}$  is the current imputation ratio.

The summations in (2.1)-(2.4) are over the  $n_r$  reported units in practice, over  $n'_r$  units in the Monte Carlo study. In addition for  $\hat{R}^{(4)}$ , the factor  $1/n'_r$  is used in the Monte Carlo study, whereas in practice the factor  $1/n_r$  is used. We also assume that the nonresponse rate is such that

$$\lim_{n'_r \rightarrow \infty} f_{n'_r} = \lim_{n'_r \rightarrow \infty} \frac{n'_r}{n_r} = f_r$$

where  $f_r$  is fixed, and  $0 < f_r < 1$ .

Lemma 1. Under above notations and assumptions, for

large  $n'_r$ ,

$$E((\hat{Y} - Y) \mid n_r) \doteq -\frac{(n'_r/n_r)}{\sum_{i=1}^{n_r} (\hat{e}_i/\pi_i)},$$

$$= 0, \text{ if } \hat{R}_{n'_r} = \hat{R}^{(1)}$$

where

$$\hat{e}_i = y_i - \hat{R}_{n_r} x_i,$$

$\hat{R}_{n_r}$  is any of the four ratio type estimators (2.1)-(2.4) using the complete data of size  $n_r$ ,  $E(\cdot \mid n_r)$  is the expectation over all possible samples of size  $n'_r$  drawn from  $n_r$ .

Proof:

Since  $\hat{y}_i = \hat{R}_{n'_r} x_i$ , it can be proved that

$\hat{R}_{n'_r} = \hat{R}_{n_r} + O_p(n'_r^{-1/2})$  for  $\hat{R}_{n'_r}$  and  $\hat{R}_{n_r}$  being any form defined in (2.1) - (2.4) using  $n'_r$  and  $n_r$  units respectively.

We then have

$$E((\hat{Y} - Y) \mid n_r)$$

$$= E\left(\sum_{i=1}^{n'_r} y_i/\pi_i + \sum_{i=1}^{n'_r} \hat{y}_i/\pi_i - \sum_{i=1}^{n'_r} y_i/\pi_i \mid n_r\right)$$

$$= n'_r E\left(\frac{1}{n'_r} \sum_{i=1}^{n'_r} (\hat{y}_i - y_i) / \pi_i \mid n_r\right)$$

$$\doteq -\frac{n'_r}{n_r} \sum_{i=1}^{n_r} \frac{1}{\pi_i} (y_i - \hat{R}_{n_r} x_i).$$

When  $\hat{R}_{n_r}$  is a form of  $\hat{R}^{(1)}$ ,

$$\sum_{i=1}^{n_r} \frac{1}{\pi_i} (y_i - \hat{R}_{n_r} x_i) = 0, \text{ and hence}$$

$$E((\hat{Y} - Y) \mid n_r) = 0.$$

Lemma 2. Under the notation and assumptions defined in this section, for large  $n'_r$ ,

$$E((\hat{Y} - Y)^2 \mid n_r)$$

$$\doteq \frac{n'_r}{n_r} \left( \sum_{i=1}^{n_r} \frac{\hat{e}_i^2}{\pi_i} + \frac{(n'_r - 1)}{(n_r - 1)} \sum_{i \neq j} \frac{\hat{e}_i}{\pi_i} \frac{\hat{e}_j}{\pi_j} \right)$$

where

$$\hat{e}_i = y_i - \hat{R}_{n_r} x_i$$

$\hat{R}_{n_r}$  is an estimator defined in (2.1) - (2.4) using the complete data set of size  $n_r$ .

Proof:

$$E((\hat{Y} - Y)^2 \mid n_r)$$

$$= E\left(\sum_{i=1}^{n'_r} (\hat{y}_i - y_i)/\pi_i \mid n_r\right)^2$$

$$= E\left(\sum_{i=1}^{n'_r} ((y_i - \hat{y}_i)/\pi_i)^2 \mid n_r\right)$$

$$+ E\left(\sum_{i \neq j}^{n'_r} ((y_i - \hat{y}_i)/\pi_i) ((y_j - \hat{y}_j)/\pi_j) \mid n_r\right)$$

$$\doteq \frac{n'_r}{n_r} \sum_{i=1}^{n_r} \frac{y_i^2}{\pi_i} - 2 \sum_{i=1}^{n_r} \frac{y_i x_i}{\pi_i} \hat{R}_{n_r} + \sum_{i=1}^{n_r} \frac{x_i^2}{\pi_i} \hat{R}_{n_r}^2$$

$$+ \frac{n'_r}{n_r} \frac{(n'_r - 1)}{(n_r - 1)} \sum_{i \neq j} \left\{ \frac{y_i y_j}{\pi_i \pi_j} - \left( \frac{y_i x_i}{\pi_i \pi_j} + \frac{y_j x_j}{\pi_i \pi_j} \right) \hat{R}_{n_r} \right.$$

$$\left. + \frac{x_i x_j}{\pi_i \pi_j} \hat{R}_{n_r}^2 \right\}$$

$$= \frac{n'_r}{n_r} \left( \sum_{i=1}^{n_r} \frac{\hat{e}_i^2}{\pi_i} + \frac{(n'_r - 1)}{(n_r - 1)} \sum_{i \neq j} \frac{\hat{e}_i}{\pi_i} \frac{\hat{e}_j}{\pi_j} \right)$$

We used the fact that

$\hat{y}_i = \hat{R}_{n'_r} x_i$ , and  $\hat{R}_{n'_r} = \hat{R}_{n_r} + O_p(n'_r^{-1/2})$  for  $\hat{R}_{n'_r}$  being any forms defined in (2.1)-(2.4). Note that if  $\hat{R}_{n_r}$  is of form of  $R^{(1)}$ , then

$$\sum_{i=1}^{n_r} \frac{\hat{e}_i^2}{\pi_i} = - \sum_{i \neq j} \frac{\hat{e}_i}{\pi_i} \frac{\hat{e}_j}{\pi_j}, \text{ because } \sum_{i=1}^{n_r} \frac{\hat{e}_i}{\pi_i} = 0.$$

Lemma 3. Under the notation and assumptions defined in this section, redefine  $\hat{y}_i = R x_i$  for  $i=1, \dots, n'_r$ , where  $R$  is a preassigned value, then

$$(1) E((\hat{Y} - Y) \mid n_r) = -\frac{n'_r}{n_r} \sum_{i=1}^{n_r} \frac{1}{\pi_i} (y_i - R x_i),$$

$$(2) E((\hat{Y} - Y) \mid n_r) = 0, \text{ iff } R = \left( \sum_{i=1}^{n_r} \frac{y_i}{\pi_i} \right) / \left( \sum_{i=1}^{n_r} \frac{x_i}{\pi_i} \right)$$

$$(3) E((\hat{Y} - Y)^2 \mid n_r)$$

$$= \frac{n'_r}{n_r} \frac{n_r}{\sum_{i=1}^{n_r} \pi_i} \frac{1}{2} (y_i - Rx_i)^2$$

$$+ \frac{n'_r}{n_r} \frac{(n'_r - 1)}{(n_r - 1)} \frac{n_r}{\sum_{i \neq j} \pi_i \pi_j} \frac{1}{2} (y_i - Rx_i)(y_j - Rx_j)$$

(4) The value of R that minimizes (3) is

$$\hat{R}_{opt} = \frac{\left( \sum_{i=1}^{n_r} \frac{x_i y_i}{\pi_i} + \frac{(n'_r - 1)}{(n_r - 1)} \sum_{i \neq j} \frac{x_i y_j}{\pi_i \pi_j} \right)}{\left( \sum_{i=1}^{n_r} \frac{x_i}{\pi_i} + \frac{(n'_r - 1)}{(n_r - 1)} \sum_{i \neq j} \frac{x_i x_j}{\pi_i \pi_j} \right)} \quad (2.5)$$

Proof: Since R is a fixed value, by the definition of  $\hat{Y}_i = Rx_i$ , and the missing at random assumption of

$n'_r$  missing units of  $y_i$ , following the similar proof in Lemma 1 and 2, we can easily derive the results.

Lemma 4. The bias and MSE of the estimated total  $\hat{Y}$  given  $n_r$ , by using  $\hat{R}_{opt}$  to impute missing  $y_i$ , is

$$E((\hat{Y} - Y) | n_r) = - \frac{n'_r}{n_r} \frac{n_r}{\sum_{i=1}^{n_r} \pi_i} \frac{1}{2} (y_i - \hat{R}_{opt} x_i) \quad ,$$

$$E((\hat{Y} - Y)^2 | n_r) = \frac{n'_r}{n_r} \frac{n_r}{\sum_{i=1}^{n_r} \pi_i} \frac{1}{2} (y_i - \hat{R}_{opt} x_i)^2$$

$$+ \frac{n'_r}{n_r} \frac{(n'_r - 1)}{(n_r - 1)} \frac{n_r}{\sum_{i \neq j} \pi_i \pi_j} \frac{1}{2} (y_i - \hat{R}_{opt} x_i) \times$$

$$(y_j - \hat{R}_{opt} x_j).$$

Proof:

$\hat{R}_{opt}$  is a function of  $n_r$  units of a complete data set

which is the population that the incomplete data samples are randomly generated in the Monte Carlo study. For a given complete data set of size

$n_r$ ,  $\hat{R}_{opt}$  is a fixed value. Following the similar proof as in Lemma 3, we have the results.

An estimator of  $\hat{R}_{opt}$  by using  $n'_r$  reported units in the Monte Carlo study is

$$\tilde{R}_{opt} = \frac{\frac{1}{n'_r} \left( \sum_{i=1}^{n'_r} \frac{x_i y_i}{\pi_i} + \frac{n'_r - 1}{n'_r - 1} \sum_{i \neq j} \frac{x_i y_j}{\pi_i \pi_j} \right)}{\frac{1}{n'_r} \left( \sum_{i=1}^{n'_r} \frac{x_i}{\pi_i} + \frac{n'_r - 1}{n'_r - 1} \sum_{i \neq j} \frac{x_i x_j}{\pi_i \pi_j} \right)}$$

$$\text{for } n'_r > 2. \quad (2.6)$$

Lemma 5. The bias and MSE of the estimated

total  $\hat{Y}$  given  $n_r$ , by using  $\tilde{R}_{opt}$  to impute missing  $y_i$  is

$$E((\hat{Y} - Y) | n_r) \doteq - \frac{n'_r}{n_r} \frac{n_r}{\sum_{i=1}^{n_r} \pi_i} \frac{1}{2} (y_i - \hat{R}_{opt} x_i)$$

$$E((\hat{Y} - Y)^2 | n_r) \doteq \frac{n'_r}{n_r} \frac{n_r}{\sum_{i=1}^{n_r} \pi_i} \frac{1}{2} (y_i - \hat{R}_{opt} x_i)^2$$

$$+ \frac{n'_r}{n_r} \frac{(n'_r - 1)}{(n_r - 1)} \frac{n_r}{\sum_{i \neq j} \pi_i \pi_j} \frac{1}{2} (y_i - \hat{R}_{opt} x_i) \times$$

$$(y_j - \hat{R}_{opt} x_j) \quad .$$

Proof:

Following the similar proofs in lemmas 1 and 2, and the fact that for large  $n'_r$ ,

$$\tilde{R}_{opt} = \hat{R}_{opt} + O_p(n_r^{-1/2}), \text{ the results follow.}$$

To use  $\tilde{R}_{opt}$  we need to know the number of nonresponse items  $n'_r$ , and the number of response items  $n_r$  in the sample. If the factor

$$(n'_r - 1)(n_r - 1)^{-1}$$

is not used in  $\tilde{R}_{opt}$ , then

$$\tilde{R}_{opt} = \frac{\sum_{i=1}^{n'_r} \frac{x_i y_i}{\pi_i} + \sum_{i \neq j} \frac{x_i y_j}{\pi_i \pi_j}}{\sum_{i=1}^{n'_r} \frac{x_i}{\pi_i} + \sum_{i \neq j} \frac{x_i x_j}{\pi_i \pi_j}} = \frac{\sum_{i=1}^{n'_r} \frac{y_i}{\pi_i}}{\sum_{i=1}^{n'_r} \frac{x_i}{\pi_i}} \quad .$$

$\tilde{R}_{opt}$  is reduced to the current imputation ratio  $R^{(1)}$ .

Another estimator of  $\hat{R}_{opt}$  is

$$R^{(5)} = \left( \sum_{i=1}^{n'_r} \frac{x_i y_i}{\pi_i} \right) / \left( \sum_{i=1}^{n'_r} \frac{x_i^2}{\pi_i} \right) \quad (2.7)$$

It can be shown that when  $R^{(5)}$   $x_i$  is used to impute missing  $y_i$ , for large  $n'_r$ , the bias and MSE of  $\hat{Y}$  for a given complete data set  $n_r$  are given in lemma 1 and 2, where

$$\hat{R}_{n_r} = \left( \sum_{i=1}^{n_r} \frac{x_i y_i}{\pi_i} \right) / \left( \sum_{i=1}^{n_r} \frac{x_i^2}{\pi_i} \right) \quad .$$

If the inclusion probability  $\pi$  is not used in (2.7), we have

$$R^{(6)} = \left( \sum_{i=1}^{n'_r} x_i y_i \right) / \left( \sum_{i=1}^{n'_r} x_i^2 \right) \quad (2.8)$$

which is the least squares estimate of R of the ratio model with constant error variance (i.e.,  $y = R x + e$ ,  $e$  is independently identically distributed with mean zero and variance  $\sigma^2$ ). It can be shown that the bias and

MSE of  $\hat{Y}$  for a given complete data set are given in lemma 1 and 2 with

$$\hat{R}_{n_r} = \left( \sum_{i=1}^{n_r} x_i y_i \right) / \left( \sum_{i=1}^{n_r} x_i^2 \right).$$

If the ordinary regression estimator is used to impute missing item  $y_i, i = 1, \dots, n'_r$ ,

$$\hat{y}_i = \hat{\alpha}_{n'_r} + \hat{\beta}_{n'_r} x_i, \quad (2.9)$$

where

$$\hat{\alpha}_{n'_r} = \bar{y}_{n'_r} - \hat{\beta}_{n'_r} \bar{x}_{n'_r}, \quad (2.10)$$

$$\bar{y}_{n'_r} = \frac{1}{n'_r} \sum_{i=1}^{n'_r} y_i,$$

$$\bar{x}_{n'_r} = \frac{1}{n'_r} \sum_{i=1}^{n'_r} x_i,$$

$$\hat{\beta}_{n'_r} = \frac{\sum_{i=1}^{n'_r} (x_i - \bar{x}_{n'_r})(y_i - \bar{y}_{n'_r})}{\sum_{i=1}^{n'_r} (x_i - \bar{x}_{n'_r})^2}, \quad (2.11)$$

then the bias and MSE of  $\hat{Y}$  for a given complete data set  $n_r$  are given in lemma 1 and 2, with

$$\hat{e}_i = y_i - \hat{\alpha}_{n_r} - \hat{\beta}_{n_r} x_i, \quad \text{and} \quad (2.12)$$

$$\hat{\alpha}_{n_r} = \bar{y}_{n_r} - \hat{\beta}_{n_r} \bar{x}_{n_r},$$

$$\bar{y}_{n_r} = \frac{1}{n_r} \sum_{i=1}^{n_r} y_i,$$

$$\bar{x}_{n_r} = \frac{1}{n_r} \sum_{i=1}^{n_r} x_i,$$

$$\hat{\beta}_{n_r} = \frac{\sum_{i=1}^{n_r} (x_i - \bar{x}_{n_r})(y_i - \bar{y}_{n_r})}{\sum_{i=1}^{n_r} (x_i - \bar{x}_{n_r})^2}.$$

The above results can be extended to more than one imputation cell if we generate the incomplete data set from the complete data set independently for each imputation cell. Let  $\hat{Y}_k, Y_k$  be the estimated totals of the incomplete data set and the complete data set respectively from imputation cell  $k, k=1, \dots, K$ . Then

$$\hat{Y} = \sum_{k=1}^K \hat{Y}_k, \quad Y = \sum_{k=1}^K Y_k.$$

Let  $n_{r_k}$  be the sample size of the reported data of imputation cell  $k$ , and we randomly suppress  $n'_{r_k}$  units from this complete data set. Let  $n_r$  be the sample size of the complete data set from all  $K$  imputation cells.

The bias of the estimated total  $\hat{Y}$  given the complete data set is

$$\begin{aligned} E((\hat{Y} - Y) \mid n_r) &= \sum_{k=1}^K E((\hat{Y}_k - Y_k) \mid n_{r_k}) \\ &= - \sum_{k=1}^K \frac{n'_{r_k}}{n_{r_k}} \frac{n_r}{\sum_{i=1}^{n_r} k} \frac{1}{\pi_{ki}} \hat{e}_{ki}, \end{aligned}$$

where

$$\hat{e}_{ki} = y_{ki} - \hat{R}_{n_{r_k}} x_{ki}, \quad k=1, \dots, K, \quad i=1, \dots, n_{r_k}.$$

$\hat{R}_{n_{r_k}}$  is a ratio estimator of (2.1) to (2.4) and (2.6) - (2.8) using the complete data set of size  $n_{r_k}$  from each

imputation cell  $k$ . For the regression estimator defined in (2.9),  $\hat{e}_i$  is defined in (2.12) for each imputation cell. Similarly, the mean square errors of the estimated total  $\hat{Y}$  given the complete data set can be written as

$$\begin{aligned} E((\hat{Y} - Y)^2 \mid n_r) &= \sum_{k=1}^K E((\hat{Y}_k - Y_k)^2 \mid n_{r_k}) \\ &= \sum_{k=1}^K \frac{n'_{r_k}}{n_{r_k}} \left( \sum_{i=1}^{n_r} k \frac{\hat{e}_{ki}^2}{\pi_{ki}} + \frac{(n'_{r_k} - 1)}{(n_{r_k} - 1)} \sum_{i \neq j} k \frac{\hat{e}_{ki}}{\pi_{ki}} \frac{\hat{e}_{kj}}{\pi_{kj}} \right). \end{aligned}$$

Under the assumption that the data are missing at random, we have already shown that the bias and MSE of the estimated total given the complete data set using various ratio and regression imputation procedures are functions of residuals of the complete data, and the nonresponse rate of each imputation cell. To compare different ratio and regression imputation procedures defined in (2.1) - (2.4) and (2.6) - (2.9) empirically, we can thus compute these biases and MSE's using Monthly Retail Trade Survey reported data and current nonresponse rates without randomly generating all possible incomplete samples.

The Monthly Retail Trade Survey reported data of December 1982 for nine SIC's were used to compare the bias and MSE of the estimated totals of the different ratio and regression type imputation procedures. The trends are calculated from the reported data of each imputation cell by these different estimators. The trends calculated by the optimum ratio procedure and the current imputation procedure are fairly close for most SIC's. (See Huang (1986)). The bias and MSE of the estimated totals by using these different imputation procedures are tabulated in Tables 1.1 and 1.2. Algebraically, we have already shown that given the complete data set, the current imputation procedure is unbiased with respect to the estimated reported total for each

imputation cell, and so are the empirical results. The relative biases (bias/estimated reported total) of the other ratio imputation procedures are relatively small, less than 3% for most data.

The optimum ratio imputation procedure,  $\tilde{R}_{opt}$ , gave the minimum mean square error among all the ratio type imputation procedures. However, the gain in efficiency of  $\tilde{R}_{opt}$  in comparing with the current imputation procedure is at most 0.002. The current imputation procedure is fairly competitive with the optimum ratio imputation procedure and is easier to compute.

Note that all the inferences of the Monte Carlo study are restricted to the data we used. The derivations of the bias and MSE are based on the assumption that the data are missing at random. The data used for the Monte Carlo study were examined to investigate the validity of this assumption. The imputation rates by sales classes of each imputation cell were calculated. There is no apparent relationship between item nonresponse rates and sales classes. The imputation rates by regions of each imputation cell were also calculated. The imputation rates are different for different regions but there is no specific pattern.

Based on the current imputation procedure, we also used mean square error (MSE) criterion to evaluate different imputation cell definitions, e.g., to answer the question of what quantiles (median, 1/4 or 1/8 or 1/16 quantiles) should be used for the cutoff of sales size classes if sales sizes are used within each firm size (group I and II) for imputation cell definition as opposed to the current fixed cutoff. The reported data for 9 SIC's from December 1982 were used. The empirical results showed that for SIC 562 the smaller the imputation cell is, the better the MSE. However, the most drastic reduction in MSE is the cell definition using 1/4 quantiles as sales cutoffs. There was an approximate 44% reduction in MSE as compared to the MSE under the current imputation cell definition. Using 1/8 quantiles as sales cutoffs a further 6% reduction over 1/4 quantiles was observed; and using 1/16 quantiles a further 3% reduction over 1/8 quantiles was observed. Overall, the empirical results varied by SIC's. In 6 of 9 SIC's, the reductions in MSE ranged from 12% (-3%) to 59% (44%) by using 1/8 (1/4) quantiles instead of the current fixed cutoff. Most of these reductions in MSE came from group II. For SIC's 541, 551, and 5813, there was little, if any, gain in using any of the quantiles considered. (See Huang (1986)).

### III. Summary

We have evaluated the bias and MSE of the estimated totals using different ratio and regression type imputation procedures (including the currently used imputation procedure) under a Monte Carlo study for a given data set.

Under the assumption that the data are missing at random, the bias and MSE of the estimated total using different ratio type imputation procedures with respect to the estimated reported total were derived for the given reported data. An optimum ratio imputation estimator was also derived along with several variants. The bias and MSE were calculated for each of nine SIC's using December 1982 retail sales data. For the given data set, the empirical results showed that the estimated total using the current imputation procedure is unbiased and has the second

smallest MSE among all ratio type imputation procedures in the study.

Since the decrease of the MSE by using the optimum imputation procedure is trivial, and extra computation and information are needed to implement this optimum imputation procedure, we do not recommend any changes of the current ratio type imputation procedure in the Monthly Retail Trade Survey.

For the given data set, there is no apparent relationship of nonresponse rate with sales within each imputation cell for all nine SIC's.

In the current imputation cells, for some SIC's, the number of establishments in Group II dominates the number in Group I; for other SIC's, the number of establishments in Group I dominates the number in Group II. The empirical results suggested that for some of the nine SIC's included in the study, we can do better by using alternative imputation cells, i.e., use sales quantiles as cutoffs within groups as opposed to the current fixed sales cutoffs. The decrease in MSE in 6 of 9 SIC's ranges from 12% to 59% by using 1/8 quantiles. We recommend that changes in the current imputation cells definition be considered, especially where empirical studies show that a significant reduction in the MSE can be achieved by increasing the number of imputation cells. We also suggest that further similar empirical studies be carried out on recent monthly data to provide a basis for changes in cell definitions for other SIC's. This will tell us whether there is a gain in using alternative imputation cells and what quantiles to use for a given SIC in a given month.

### VI. Acknowledgement

The author gratefully acknowledges the helpful comments and criticisms of Cary Isaki, Nash Monsour and Carl Konschnik.

### V. References

1. Bailey, L., Chapman, D.W. and Kasprzyk, D. (1985). "Nonresponse Adjustment Procedures at the Bureau of the Census: A Review. Paper presented at First Annual Census Bureau Research Conference, Reston, Virginia.
2. Cassel, C.M., Sarndal, C.E., and Wretman, J.H. (1979). "Some Uses of Statistical Models in Connection with the Nonresponse Problem," Incomplete Data in Sample Surveys. Volume 3. W.G. Madow, I. Olkin and D.B. Rubin, eds., Academic Press: New York pg. 143-160.
3. Cochran, W.G. (1977). Sampling Techniques. John Wiley and Sons, New York.
4. Fuller, W.A. (1976). Introduction to Statistical Time Series. John Wiley and Sons, New York.
5. Huang, E.T. (1984). "An Imputation Study for the Monthly Retail Trade Survey," Proceedings of the American Statistical Association, Section on Survey Research Methods, pp. 610-615.
6. Huang, E.T. (1986). "Report on the Imputation Research for the Monthly Retail Trade Survey" Statistical Research Division report series, report number: CENSUS/SRD/RR-86-09/. Bureau of the Census.
7. Little, R.J.A (1986). "Missing Data in Census Bureau Surveys". Paper presented at Second Annual Census Bureau Research Conference, Reston, Virginia.

TABLE 1.1 The Bias (Relative Bias (%)) of the Estimated Total  
By Using Different Imputation Procedures  
December 1982

Unit: U.S. Dollars

SIC	n	Estimated Reported Total	R(1)	R(2)	R(3)	R(4)	$\tilde{R}_{opt}$	R(5)	R(6)	Regression Estimator
562 (Women's Ready-to-Wear Stores) (%)	1445	1,636,658,834	0 (0)	4,745,104 (0.290)	18,303,757 (1.118)	29,244,857 (1.787)	-495,262 (-0.030)	-18,584,772 (-1.136)	-49,527,438 (-3.026)	73,566,518 (4.495)
521 (Building Materials Stores) (%)	635	1,933,849,833	0 (0)	7,282,673 (0.377)	14,289,364 (0.739)	24,869,789 (1.286)	-133,689 (-0.007)	-5,946,523 (-0.307)	-850,347 (-0.044)	8,085,767 (0.418)
531 (Department Stores) (%)	7557	14,758,285,090	0 (0)	202,901 (0.001)	71,895,263 (0.487)	72,247,748 (0.490)	-340,476 (-0.002)	-77,969,143 (-0.528)	-78,411,093 (-0.531)	425,515 (0.003)
541 (Grocery Stores) (%)	2428	12,374,995,572	0 (0)	51,782,545 (0.418)	-13,933,939 (-0.113)	24,927,603 (0.201)	-409,814 (-0.003)	-1,945,308 (-0.016)	59,579,151 (2.277)	39,159,736 (0.316)
551 (Motor Vehicle Dealers) (%)	753	14,565,413,603	0 (0)	14,544,341 (0.100)	59,652,581 (0.410)	61,644,105 (0.423)	-3,169,391 (-0.022)	-53,164,307 (-0.365)	-5,047,197 (-0.035)	57,052,376 (0.392)
572 (Household Appliances, Radio, TV Stores) (%)	500	571,806,693	0 (0)	6,876,815 (1.203)	2,920,688 (0.511)	5,941,472 (1.039)	-266,227 (-0.047)	-4,071,594 (-0.712)	13,020,508 (2.277)	-2,565,749 (-0.449)
5812 (Eating Places) (%)	1531	6,055,819,018	0 (0)	-2,022,721 (-0.033)	31,094,652 (0.513)	47,959,481 (0.792)	-720,701 (-0.012)	-18,773,433 (-0.310)	-41,890,604 (-0.692)	24,760,883 (0.409)
5813 (Drinking Places) (%)	420	642,146,909	0 (0)	-163,904 (-0.026)	-647,688 (-0.101)	-151,990 (-0.024)	59,308 (0.009)	1,057,819 (0.165)	634,619 (0.099)	-2,087,831 (-0.325)
592 (Liquor Stores) (%)	542	1,740,095,873	0 (0)	-3,303,540 (-0.190)	3,987,841 (0.229)	1,672,854 (0.096)	-345,361 (-0.020)	-2,725,450 (-0.157)	-18,100,886 (-1.040)	2,683,752 (0.154)

TABLE 1.2 The MSE of the Estimated Total By Using Different Imputation Procedures  
(And the Ratio to its Current Imputation Procedure)  
December 1982

Unit: \$10<sup>6</sup>

SIC	n	R(1)	R(2)	R(3)	R(4)	$\tilde{R}_{opt}$	R(5)	R(6)	Regression Estimator
562 (Women's Ready-to-Wear Stores)	1445	122,250,188 (1)	149,423,484 (1.2223)	485,083,727 (3.9680)	542,620,951 (4.4386)	122,151,733 (0.9992)	254,714,015 (2.0835)	1,384,741,428 (11.327)	2,236,292,027 (18.293)
521 (Building Materials Stores)	635	373,293,260 (1)	508,178,456 (1.3613)	634,679,773 (1.7002)	810,125,338 (2.1702)	373,238,170 (0.9999)	405,332,341 (1.0858)	425,909,403 (1.1410)	482,528,485 (1.2926)
531 (Department Stores)	7557	148,129,951 (1)	148,187,203 (1.0004)	5,362,890,768 (36.204)	5,363,246,499 (36.206)	148,069,922 (0.9996)	6,176,598,292 (41.697)	6,178,473,892 (41.710)	121,809,138 (0.8223)
541 (Grocery Stores)	2428	975,155,627 (1)	2,732,057,367 (2.8017)	1,383,526,763 (1.4188)	2,492,746,633 (2.5563)	975,070,021 (0.9999)	1,306,876,310 (1.3402)	2,902,580,498 (2.9765)	2,341,096,627 (2.4007)
551 (Motor Vehicle Dealers)	753	2,617,510,636 (1)	2,906,824,784 (1.1105)	4,919,457,310 (1.8794)	5,040,731,723 (1.9258)	2,610,996,151 (0.9975)	4,108,583,221 (1.5697)	2,689,858,236 (1.0276)	5,023,064,502 (1.9190)
572 (Household Appliances, Radio/TV Stores)	500	25,604,966 (1)	60,624,441 (2.3677)	41,244,917 (1.6108)	52,676,300 (2.0573)	25,562,855 (0.9984)	35,419,411 (1.3833)	134,778,334 (5.2638)	83,056,483 (3.2438)
5812 (Eating Places)	1531	410,794,971 (1)	477,812,049 (1.1631)	1,134,373,917 (2.7614)	1,339,936,313 (3.2618)	410,512,304 (0.9993)	551,848,720 (1.3434)	1,168,011,026 (2.8433)	669,950,846 (1.6309)
5813 (Drinking Places)	420	4,150,594 (1)	4,205,883 (1.0133)	4,691,465 (1.1303)	4,309,152 (1.0382)	4,145,352 (0.9987)	5,229,200 (1.2599)	4,577,926 (1.1030)	8,792,564 (2.1184)
592 (Liquor Stores)	542	110,351,071 (1)	150,489,266 (1.3637)	183,571,784 (1.6635)	118,555,857 (1.0744)	110,152,107 (0.9982)	119,873,690 (1.0863)	428,523,608 (3.8833)	144,471,410 (1.3092)