

The Effect of Content Errors on Bias and Nonsampling Variance in Estimates Derived From Samples of Administrative Records

Barry W. Johnson and Darien B. Jacobson, Internal Revenue Service

The Statistics of Income Division (SOI) of the Internal Revenue Service (IRS) uses a number of methods for ensuring the quality and integrity of the data it produces for tax administration research. As a first line of quality assurance, codes and mathematically related data items are extensively tested as SOI employees enter them into computer databases. In addition, for a subsample of returns selected and processed in most studies, SOI assigns a second employee to reenter and edit the data. Values from the first and second edit are then computer-matched. A supervisor resolves discrepancies discovered during the match. The original value, second value, and correct values are all collected as a part of the quality review system, as are a set of codes that describe the cause of the error, in broad categories.

This paper will use quality review data from Federal estate tax returns (Form 706) selected into the Calendar Year 2002 SOI Estate Tax Study to estimate the effects of nonsampling error on estimates derived from the final data file.

► Background

The Federal estate tax is levied on estates for the right to transfer assets from a decedent's estate to its beneficiaries; it is not an inheritance tax. A Federal estate tax return must be filed for every U.S. decedent whose gross estate, valued on the date of death, combined with certain lifetime gifts made by the decedent, equals or exceeds the filing threshold applicable for the decedent's year of death. A decedent's estate must file a return within 9 months of a decedent's death, but a 6-month extension is usually granted.

All of a decedent's assets, as well as the decedent's share of jointly owned and community property assets, are included in the gross estate for tax purposes and reported on Form 706. Also reported are most life insurance proceeds, property over which the decedent possessed a general power of appointment, and certain transfers made during life.

Expenses and losses incurred in the administration of the estate, funeral costs, and the decedent's debts are allowed as deductions against the estate for the purpose of calculating the tax liability. A deduction is allowed for the full value of bequests to the surviving spouse. Bequests to qualified charities are also fully deductible.

► Data Description

The 2002 SOI Estate Tax Study was a stratified, random sample of returns filed in Calendar Year 2002 and was the second year in a 3-year study of Federal estate tax returns filed 2001-2003. The sample was designed for use in both estimating tax revenues in all 3 calendar years and personal wealth holdings for 2001 decedents. The 3-year sample period was devised to ensure that nearly all returns filed for 2001 decedents would be subjected to sampling, since a return could be filed up to 15 months after the decedent's death. The design had three stratification variables: size of total gross estate plus the value of most taxable gifts made during the decedent's life, age at death, and year of death. The year-of-death variable was separated into two categories, 2001 year of death and non-2001 year of death, in order to facilitate studies of 2001 decedents. Returns were chosen before audit examination and selected using a stratified random probability sampling method. A portion of the sample was selected because the ending digits of the decedents' Social Security Numbers (SSN) corresponded with those in the 1-percent Social Security Administration Continuous Work History Sample. However, the majority of returns were selected on a flow basis using the Bernoulli sampling method.

The sampling mechanism was a permanent random number based on an encryption of the decedent's SSN. Sample rates were preset based on the desired sample size and an estimate of the population. Sampling rates ranged from 3 percent to 100 percent, with more than half of the strata selected with certainty.

Data collection for the 2002 Estate Tax Study was

conducted at the IRS Cincinnati Submission Processing Center. Employees entered the data from the estate tax return into a database using a Graphical User Interface (GUI) data entry system. Nearly 100 distinct data items were captured, with some balance sheet assets recurring hundreds, even thousands, of times, as assets were allocated to 32 different categories, such as stocks, bonds, and real estate. Tax returns ranged in size from a dozen to many thousands of pages, including appraisals, investment account listings, and legal documents. Tests embedded in the data entry system were used to validate entries and to ensure that mathematical relationships among variables were correctly preserved. There were more than 200 validation tests performed on each tax return included in the 2002 study.

While embedded testing can assure that codes are correct within a given range of values and that fields are mathematically consistent, many of the decisions that employees make when transforming tax return information into statistically usable data are not easily tested. For example, while several codes may be valid, determining the best code to describe a particular taxpayer's behavior or characteristics cannot always be automated. To address this problem, SOI developed a double entry quality review system. This system is a valuable tool for measuring both individual employee performance and overall data quality.

► **Quality Review System**

A subsample of returns in the 2002 Estate Tax Study was subjected to additional review for quality assurance purposes. Returns were included in the quality review (QR) subsample through two different mechanisms, 100-percent review and product review. The 100-percent review consisted of all returns that were edited while an employee was in training. Product review was selected after the training period had been completed, and it comprised a 10-percent random sample of each employee's work. The product review sample was selected on a flow basis method using a pseudorandom number called the Transform Taxpayer Identification Number, or TTIN. The TTIN is a unique random number that is generated by mathematically transforming selected digits of the decedent's Social Security Number. The TTIN was then compared to the sample number, which represented the

sample rate, in this case, 10 percent. If the TTIN was less than the sample number, then the return was selected for product review.

Under the double-entry quality review system, one return was entered into the computer system twice by two different employees. The first employee did not know that a return was selected for review until after the first edit was complete, and the second employee was not allowed to see the first employee's entries. Therefore, each return had two versions in the database, the first edit and the second edit, and each was entered independently of the other.

When both employees finished editing a return, the computer compared the values from the original and QR versions. In some cases, the two versions matched perfectly; so, the return was released from the system, and the first edit data was treated as final and stored for later analysis. However, if mismatches between the two versions occurred, the discrepancies were stored in a separate data table to be reviewed by a supervisor.

The supervisor reviewed the discrepancies and charged the errors, assigning two codes to each discrepancy--one to identify the incorrect value and the other to describe the cause of the error. A discrepancy code was assigned to the error to explain which version was considered incorrect. Discrepancy codes were assigned to one of the following: the first version, the second version, both versions, or neither version. An error was assigned to both versions if both of the employees entered or interpreted the information from the return incorrectly. In this case, the supervisor was also required to supply the correct data value. In some cases an error was not assigned to either version, usually when the discrepancy was the result of a data processing peculiarity and not a true database error. After the error was assigned a discrepancy code, a numeric error resolution code was assigned to describe why the entry was incorrect. Error resolution codes indicate situations such as spelling errors, incorrect money amounts, or incorrectly assigned codes.

Once the supervisor reviewed all the discrepancies, each employee was given a list of the discrepancies, along with the discrepancy and error resolution codes,

so that any first edit errors detected during quality review could be corrected prior to considering return processing complete. The feedback from the review also enabled employees to learn from their mistakes on each return and carry this knowledge into the editing of other returns. In the end, there is a database consisting of a table that includes all the values from the second edit of the return as entered, a quality review table containing a record of each discrepancy between the first and second edits (along with codes indicating who made the error and why), and a final data table containing the correct version of the return data that will ultimately be sent to customers.

For this paper, only a portion of the quality review data was used for analysis. First, data that were collected during periods of training, 100 percent review, were excluded. Second, only errors that were charged to the first edit or to both edits, meaning that the error required a correction to the final data set, were retained. This was done because these errors are more representative of errors that remain in the roughly 90 percent of the 2002 estate tax sample that was not selected for quality review. Third, errors that reflected idiosyncrasies related to the edit process itself, and not true data errors, were eliminated.

► **Empirical Results**

Quarterly accuracy rates for each employee who worked on the Estate Tax Study for 2002 were generated using the product review data (see Figure 1). These rates were calculated using the number of returns that had at least one error charged to the first edit divided by the total number of returns that had been selected for quality review. The accuracy rates for all of the employees are not very high. However, these rates are a return level measure; any return with one or more errors is considered incorrect. The Form 706 includes an average of 150 data entry fields, while complex returns can have more than a thousand entries; so, the probability of making just one mistake is very high. In fact, the average number of errors for each return is only 6.3.

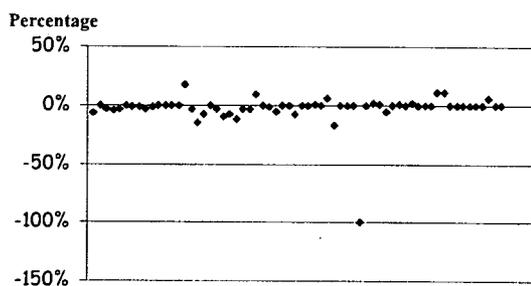
Traditionally, supervisors have focused quality improvement efforts on those fields that are in error most frequently. By looking at the occurrence of variables *ex-*

Figure 1: Employee Accuracy Rates

Employee	Accuracy Rates			
	Quarter 1	Quarter 2	Quarter 3	Quarter 4
17000	46.3%	23.9%	41.7%	21.7%
17100	25.0%	0.0%	0.0%	0.0%
17200	29.2%	30.8%	31.9%	40.0%
17300	57.1%	100.0%	91.7%	33.3%
17400	52.1%	28.6%	50.0%	37.9%
17500	44.4%	24.1%	54.8%	0.0%
17600	42.2%	51.9%	33.9%	46.2%
17700	41.9%	28.6%	39.3%	34.5%
17800	49.1%	25.0%	58.5%	45.6%
17900	52.3%	34.3%	59.0%	50.0%
17001	23.1%	34.2%	18.6%	44.7%
17002	39.2%	33.3%	36.2%	45.0%
17003	22.9%	20.7%	37.8%	29.1%
17004	34.2%	31.6%	22.0%	72.7%
17005	30.8%	0.0%	0.0%	37.9%
17006	26.5%	27.7%	41.4%	42.9%

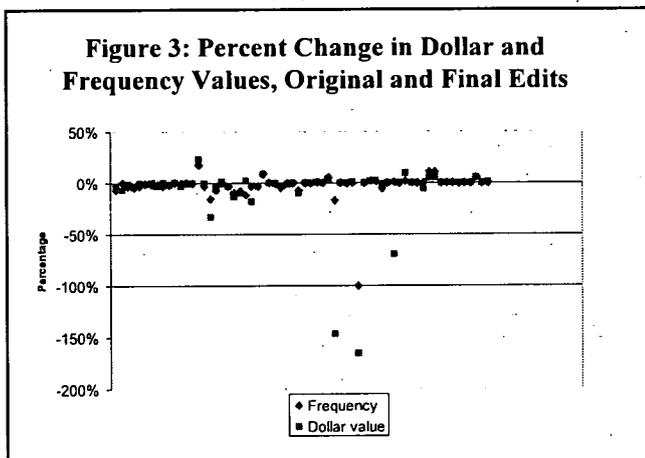
ante, using the first edit data, and *ex-post*, using the final corrected data file, it is possible to identify the frequency of original edit errors in the quality review sample. Figure 2 shows the percent changes in frequencies for variables on the file; each diamond represents a different variable. Frequencies change because many variables on the file represent balance sheet items, assets like stocks, bonds, mutual funds, and various types of real estate, which are not necessarily present in each decedent's portfolio. When an asset is incorrectly classified, not only does it change the dollar value of estimate, it also changes the frequency of occurrence of that particular attribute or asset type in the population estimates. This can be particularly problematic if the asset is of special interest to researchers. For example, there has been much discussion in the press about providing estate tax relief to small business owners. Errors that either under- or overcount the number of estates that have small

Figure 2: Percent Change in Frequencies, Original and Final Edits



businesses could have an impact on this debate. The percentages shown on the graph represent the aggregate correct frequency in the overall quality review sample, less the aggregate number originally reported, divided by the correct number. Negative percentages indicate cases where an asset was incorrectly included on the first edit. For example, the first employee may have incorrectly classified a balance sheet entry as a publicly traded stock, while the second employee may have correctly classified it as a mutual fund invested in a mix of financial assets. The percent changes in frequencies are generally close to zero, but there are some notable outliers.

Figure 3 shows percentage changes in dollar amounts between first and second edits overlaid on the frequency differences shown in Figure 2. Each point represents a single variable on the file. While the pattern for the dollar differences is similar to that of the frequencies, with many differences close to zero, the magnitude of the dollar differences is larger for several variables. There are two variables for which the original entries resulted in aggregate dollar values that were overstated by roughly 150 percent. This highlights the potentially large effects on final estimates that can arise from even one large dollar value error, especially for variables that are not widely distributed in the overall population. Thus, it is important to monitor both the size and frequency of data entry errors.



Unweighted error statistics are clearly useful for monitoring data quality and assessing opportunities for operational improvements during a study period. However, since the SOI study of Federal estate tax returns is based on a stratified random sample of the filing popu-

lation, the effect of data entry error on final population estimates derived from this sample will vary inversely with the selection rate associated with each return. Using appropriate sample weights, it is possible to use the 10-percent QR sample to estimate the effects of data entry errors on population estimates derived from the remaining 90 percent of the returns in the final SOI data file that were not subjected to double-entry quality review. Weighted estimates provide a different perspective on the effects of nonsampling error due to the nature of the underlying estate study sample and the fact that the financial characteristics of estate tax decedents vary greatly among age and wealth classes. For example, younger decedents and those with large estates are selected into the estate tax sample with certainty and comprise more than 40 percent of the total sample file. Both groups of decedents are more likely to have had portfolios that are more complex and, thus, more subject to data entry errors than their either less wealthy, or older, cohorts. This is because many older wealth holders convert their portfolios to assets that produce tax-preferred income, usually resulting in returns that contain fewer business arrangements, which are more difficult to classify than market assets. Because the quality review sample is not stratified, weighted estimates will provide a more balanced measure of the overall effects of data entry errors on final estimates. Weighted estimates for the quality review sample were generated by using the design-based weight from the stratified estate study sample (W_s), multiplied by a quality review weight (W_q). The quality review weight itself was developed by first post-stratifying the quality review samples within the original selection strata as indicated below [1]:

$$\text{Final Weight} = W_s * W_q$$

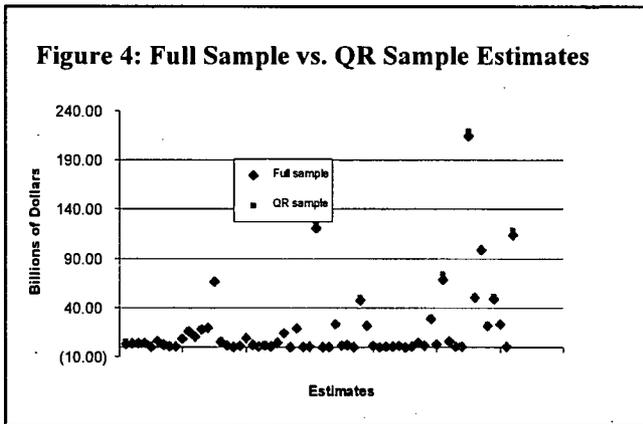
$$\text{Where } W_s = N_i/n_i$$

$$\text{Post-Stratification: } W_q = n_{if}/n_{qif}$$

For some strata, the quality review sample was either zero or too small to create a post-strata cell. For these cases, strata were collapsed across age categories so that estate size classes were preserved.

Figure 4 shows full population dollar value estimates from the quality review data using the post-stratified quality review weight and compares them to population

estimates using the full weighted estate study sample. Each pair of data points represents a different variable on the file. The quality review data estimates for each variable are denoted by the gray squares, and the full sample estimates are denoted by the black diamonds. For most variables, the QR sample estimates are larger than the population estimates from the full estate sample, indicating that the QR sample introduces a positive bias. This bias arises because the QR sample is a simple random sample of a stratified sample that favors large dollar value returns. In such cases, ratio raking can often be employed to decrease the bias; however, in this case, the QR sample size was insufficient in the lower gross estate size classes.



While the weighted QR data estimates are somewhat biased due to the design of the sample, they still provide an important indication of the effects of data entry errors on final estate tax sample estimates. Figure 5 shows weighted and unweighted estimates of aggregate differences between original and final values of both frequency and dollar value estimates for selected variables. A negative value means that a variable was overrepresented in the original, uncorrected data, and a positive value means it was originally underrepresented. Weighted results rank errors differently for some of the variables. For example, errors in classifying noncorporate business assets had a much greater impact on final weighted estimates than would have been evident had the analysis been limited to examining the unweighted QR data. Conversely, the unweighted QR data implied that the effects of errors on estimates of farm real estate were greater than they are in the final, weighted estimates. Clearly, using weighted estimates, along with

the unweighted quality review data, provides a more balanced method of assessing where to focus data quality improvement efforts.

Figure 5: Differences Between First and Final Edit

Data Element	Frequency	Dollar Value
Non-corporate businesses	-11.00%	-5.79%
	<i>-5.29%</i>	<i>-3.55%</i>
Publicly traded stock	15.02%	20.00%
	<i>15.38%</i>	<i>23.40%</i>
Closely held stock	-3.06%	-1.01%
	<i>-3.42%</i>	<i>-.71%</i>
Real estate	6.70%	7.34%
	<i>6.82%</i>	<i>6.17%</i>
Farm land	-.91%	-1.09%
	<i>-1.95%</i>	<i>-3.66%</i>
Funeral expenses	.25%	.15%
	<i>.09%</i>	<i>.04%</i>

Values in *italics* are unweighted

Figure 6 compares the weighted percent differences between original edit estimates and final, corrected estimates with coefficients of variation (C.V.) from the full estate tax study sample in order to relate the sampling and nonsampling variances associated with selected fields. For some estimates, such as the values for non-corporate businesses and publicly traded corporations, the nonsampling error attributable to data entry is much greater than the sampling variance. For others, such as estimates of stock in closely held or untraded corporations and farm land, the sampling error, represented by the C.V., is actually greater than the nonsampling error attributable to data entry errors, indicating that data entry errors are not a significant cause of additional variance in the estimates. Fields for which nonsampling error

Figure 6: Data Entry Error vs. Sample Variance

Data Element	Frequency		Money Amount	
	% diff	C.V.	% diff	C.V.
Non-corporate businesses	-11.00%	4.45%	-5.79%	3.89%
Publicly traded stock	15.02%	.78%	20.00%	1.17%
Closely held stock	-3.06%	3.47%	-1.01%	2.18%
Real estate	6.70%	1.92%	7.34%	2.19%
Farm land	-.91%	4.34%	-1.09%	4.68%
Funeral expenses	.25%	.57%	.15%	1.19%
Spousal trusts	4.25%	2.97%	1.29%	1.58%

is relatively large provide opportunities for future data quality improvement efforts.

► **Conclusion**

There is much to be learned through careful analysis of the data generated by SOI's double-entry quality review systems. The results of these analyses can be used to improve data collection systems and enhance worker training. Information on nonsampling error should also be useful to data users who could use data quality metrics to more accurately interpret economic modeling results and to ultimately build models that are more robust.

This analysis, however, revealed that the database format and the type of data that are collected from the quality review samples make certain types of analysis difficult, if not impossible. While a complete copy of the second edit is saved for all QR returns, the original, uncorrected first edit values are not saved when first edit errors require corrections. Information on discrepancies is kept in all cases, but, because corrections can involve changing any number of related fields, it is difficult to reconstruct exactly the first employee's original entries. If more sophisticated analysis is desired, including the study of secondary errors that arise as a result of a primary data entry error, archiving a complete copy of the first edit, along with associated error reason and discrepancy codes, should be considered.

It is also important that supervisors apply error reason and discrepancy codes consistently. All too often, discrepancies are resolved by several different supervisors. Some, especially those serving in a temporary capacity, may feel a great deal of peer pressure to avoid assigning errors to individual employees, even in cases where the assignment of an error would not directly impact employee performance appraisals, such as when an error is attributable to lack of clarity in editing instructions. This inconsistency makes it difficult to measure

the extent to which errors exist and to learn of ways to avoid them in the future.

Related to this problem is that the measure of employee performance currently in place is not adequate. It is simply unfair to use a return level measure of accuracy when the difficulty of the work is so variable across returns. A more balanced measure would relate the number of individual errors an employee makes to the number of fields he or she actually edited, thus giving full consideration to the number of edit decisions that were made on each return.

Finally, there are sample design issues that became apparent from this analysis. The QR sample is biased and could be improved by taking into consideration the underlying structure of the estate tax study sample design. Even this would not provide coverage of variables that are relatively rare, but perhaps important, in policy debates. To address this problem, samples could either be increased or targeted to include more returns with important characteristics, such as those filed for small business owners, or returns that, because of the types of entries made during first edit, are more likely to contain significant problems. Samples could also vary with worker skill levels. One possibility would be to develop a system that sets a weekly QR sample rate for each individual employee based on individual rolling average accuracy rates. Sample rates could be set automatically based on preset performance standards. Automating the process would avoid putting supervisors in the awkward position of having to 'punish' poor performers with additional oversight, making it easier to match feedback and training efforts to performance levels.

► **Footnote**

- [1] The subscript "if" signifies that certain reject returns were removed from the estate study sample prior to post-stratifying.