

Section 3

Description of the Sample and Limitations of the Data

This section describes the 1996 Corporate sample design, including the methods used in the selection of returns, data capture, data cleaning, and data completion. The techniques used to produce estimates and an assessment of the data limitations, including measures of sampling variability, are also discussed.

Background

From Tax Year 1916 through Tax Year 1950, data were extracted for the Statistics of Income (SOI) program from each corporate return filed. Stratified probability sampling was introduced for Tax Year 1951. Since then, the size of the samples has generally decreased while the population has increased. For example, for Tax Year 1951 the sample comprised 41.5 percent of the entire population, or 285,000 of the 687,000 total returns filed. For 1996, the sample proportion had decreased to about 1.9 percent of the total population of over 4.9 million.

For 1951, stratification was by size of total assets and industry. From 1952 through 1967, the stratification was by size only. The size was measured by volume of business (1953-1958) or total assets (1952, and 1959-1967). Since 1968, returns have been stratified by both total assets and a measure of income, the definition of which depends on the return's form type [1].

Target Population

The target population consists of all returns of active corporations organized for profit that are required to file one of the 1120 forms that are part of the SOI study.

Survey Population

The survey population includes the returns that filed on one of the 1120 forms in the SOI study and posted to the IRS Business Master File (BMF). Amended returns and returns that changed because of a tax audit are excluded. The following table gives the actual number of corporate returns by form type that were subject to sampling during Tax Years 1993 through 1996. These population counts will differ from all the estimated population counts in this publication because they include out-of-scope returns that are excluded from the tabulations (see page 10).

Form Type	Tax Year			
	1993	1994	1995	1996
1120	1,980,483	2,214,657	2,235,287	2,232,069
1120-A	325,773	321,402	325,249	289,477
1120S	2,011,167	2,139,353	2,267,178	2,420,886
1120-L	1,942	1,829	1,718	1,636
1120-PC	2,760	2,846	2,928	3,124
1120-RIC	6,931	7,712	8,478	8,731
1120-REIT	354	394	473	534
1120-F	11,274	11,905	10,875	11,879
Total	4,340,684	4,700,098	4,852,186	4,968,336

Sample Design

The current sample design is a stratified probability sample, with stratification by form type, and either size of total assets alone, or both size of total assets and a measure of income. Forms 1120 and 1120-A are stratified by size of total assets and size of "proceeds." Size of "proceeds" is used as the measure of income, and is defined to be the larger of the absolute value of net income (or deficit) or the absolute value of "cash flow," which is the sum of net income and several depreciation amounts. Forms 1120-F, 1120-L, 1120-PC, 1120-RIC, and 1120-REIT are each stratified by size of total assets only. Form 1120S is stratified by size of total assets and, as the measure of income, size of ordinary income.

The design process began with projected population totals derived from those used to estimate IRS administrative workloads and are adjusted based on previous years' population distributions. Using projected population totals by sampling strata, an optimal allocation, based on variance and cost estimates, was carried out to assign sample rates such that the overall projected sample size is 92,000. A Bernoulli sample is selected independently from each stratum with rates ranging from .25 percent to 100 percent. The total realized sample for Tax Year 1996, including inactive corporations and rejected returns, is 94,325 returns. Figure C on the following page shows the stratum boundaries, sampling rates, and population and sample counts for each form type. The table also shows the adjusted population and sample counts after reclassification of returns due to errors in the stratifying variables (see subsection on Processing Errors, page 13, for further information on the handling of mis-stratified returns).

Bertrand Überall, Richard Collins, and Valerie Puckett were responsible for the sample design and estimation of the SOI 1996 Corporation Program under the direction of Yahia Ahmed, Chief, Mathematical Statistics Section, Statistical Computing Branch.

1996 Corporation Returns - Sample Description and Data Limitations

Figure C. -- Corporation Returns: Number Filed, Number in Sample, and Sampling Rates by Sample Selection Class

Sample Class Number	Description of Sample Selection Classes		Sampling Rates (%)	Number of Returns			
				BMF counts		After Adjustments**	
	Size of Total Assets	Size of Proceeds*		Population***	Sample	Population****	Sample****
	All Returns, Total			4,968,336	94,325	4,968,490	94,172
	Form 1120 w/ Form 5735 attached, Total			434	434	434	433
1	Under \$100,000,000		100.00	349	349	349	348
2	\$100,000,000 - \$250,000,000		100.00	39	39	38	38
3	\$250,000,000 or more		100.00	46	46	47	47
	Form 1120 (no Form 5735 attached), 1120-A, Total			2,521,112	58,436	2,521,216	58,339
4	Under \$50,000	Under \$25,000	0.30	951,521	2,877	948,786	2,899
5	\$50,000 - \$100,000	\$25,000 - \$50,000	0.37	351,511	1,305	351,828	1,341
6	\$100,000 - \$250,000	\$50,000 - \$100,000	0.57	449,667	2,525	452,675	2,591
7	\$250,000 - \$500,000	\$100,000 - \$250,000	1.20	290,694	3,461	291,375	3,502
8	\$500,000 - \$1,000,000	\$250,000 - \$500,000	2.00	195,927	3,938	195,815	4,001
9	\$1,000,000 - \$2,500,000	\$500,000 - \$1,000,000	4.60	145,031	6,662	144,816	6,726
10	\$2,500,000 - \$5,000,000	\$1,000,000 - \$1,500,000	6.00	55,093	3,331	54,987	3,357
11	\$5,000,000 - \$10,000,000	\$1,500,000 - \$2,500,000	11.00	30,301	3,177	30,088	3,192
12	\$10,000,000 - \$25,000,000	\$2,500,000 - \$5,000,000	30.00	21,303	6,327	21,046	6,257
13	\$25,000,000 - \$50,000,000	\$5,000,000 - \$10,000,000	50.00	10,540	5,309	10,457	5,233
14	\$50,000,000 - \$100,000,000	\$10,000,000 - \$15,000,000	100.00	6,885	6,885	6,815	6,745
15	\$100,000,000 - \$250,000,000	\$15,000,000 or more	100.00	6,779	6,779	6,697	6,664
16	\$250,000,000 - \$500,000,000		100.00	2,313	2,313	2,293	2,293
17	\$500,000,000 or more		100.00	3,547	3,547	3,538	3,538
	Form 1120S, Total			2,420,886	24,864	2,420,922	24,854
18	Under \$50,000	Under \$25,000	0.25	1,027,358	2,464	1,018,433	2,502
19	\$50,000 - \$100,000	\$25,000 - \$50,000	0.30	383,315	1,160	385,263	1,208
20	\$100,000 - \$250,000	\$50,000 - \$100,000	0.45	413,434	1,859	417,598	1,928
21	\$250,000 - \$500,000	\$100,000 - \$250,000	1.00	258,066	2,625	258,733	2,654
22	\$500,000 - \$1,000,000	\$250,000 - \$500,000	1.60	147,767	2,357	148,699	2,387
23	\$1,000,000 - \$2,500,000	\$500,000 - \$1,000,000	3.70	106,921	3,946	107,714	3,987
24	\$2,500,000 - \$5,000,000	\$1,000,000 - \$1,500,000	4.60	42,672	2,002	43,222	2,022
25	\$5,000,000 - \$10,000,000	\$1,500,000 - \$2,500,000	9.00	23,345	2,013	23,638	2,030
26	\$10,000,000 - \$25,000,000	\$2,500,000 - \$5,000,000	24.00	12,579	3,026	12,512	2,991
27	\$25,000,000 - \$50,000,000	\$5,000,000 - \$10,000,000	40.00	3,362	1,345	3,259	1,296
28	\$50,000,000 - \$100,000,000	\$10,000,000 - \$15,000,000	100.00	1,131	1,131	1,075	1,073
29	\$100,000,000 - \$250,000,000	\$15,000,000 or more	100.00	719	719	638	638
30	\$250,000,000 or more		100.00	217	217	138	138
	Form 1120-L, Total			1,636	905	1,639	901
31	Under \$50,000,000		43.00	1,290	559	1,263	525
32	\$50,000,000 - \$250,000,000		100.00	118	118	117	117
33	\$250,000,000 or more		100.00	228	228	259	259
	Form 1120-F (with effectively-connected income in U.S.), Total			11,879	1,974	11,882	1,973
34	Under \$50,000,000		14.00	11,516	1,611	11,514	1,605
35	\$50,000,000 - \$100,000,000		100.00	91	91	93	93
36	\$100,000,000 or more		100.00	272	272	275	275
	Form 1120-PC, Total			3,124	1,206	3,126	1,186
37	Under \$50,000,000		30.00	2,780	862	2,784	844
38	\$50,000,000 - \$250,000,000		100.00	212	212	210	210
39	\$250,000,000 or more		100.00	132	132	132	132
	Form 1120-REIT, Total			534	435	534	427
40	Under \$50,000,000		50.00	213	114	221	114
41	\$50,000,000 - \$250,000,000		100.00	133	133	133	133
42	\$250,000,000 or more		100.00	188	188	180	180
	Form 1120-RIC, Total			8,731	6,071	8,737	6,059
43	Under \$50,000,000		25.00	3,542	882	3,536	872
44	\$50,000,000 - \$100,000,000		100.00	1,133	1,133	1,138	1,132
45	\$100,000,000 - \$250,000,000		100.00	1,632	1,632	1,643	1,635
46	\$250,000,000 - \$500,000,000		100.00	939	939	934	934
47	\$500,000,000 or more		100.00	1,485	1,485	1,486	1,486

* Proceeds is defined as the larger of absolute value of net income (deficit) or absolute value of cash flow (depreciation + depletion + net income).

** These adjustments include restratification (see section on Processing Errors, page13).

*** Includes added returns not posted to the BMF during the two-year IRS processing period.

**** Does not include missing returns, but does include added returns not posted to the BMF during the two year IRS processing period.

Note: Returns were classified according to either size of total assets or size of proceeds, whichever corresponded to the higher sample class.

EXAMPLE: A Form 1120 return with total assets of \$750,000 and a proceeds of \$75,000 is in sample class 8 (based on total assets), rather than in sample class 6 (based on proceeds).

Sample Selection

Corporation income tax returns are filed at the ten IRS service centers located throughout the country. All corporate returns are processed initially to determine tax liability and are then made available for other programs including SOI. All tax data are transmitted and updated on a weekly basis to the IRS Business Master File (BMF) system located in Martinsburg, West Virginia. This system serves as the point of selection for the sample, which was selected on a weekly basis.

Sample selections for Tax Year 1996 occurred over the period of July 1996 through June 1998. A 24-month sampling period is needed for two reasons. First, approximately 22.7 percent of all corporations have noncalendar year accounting periods. In order to take the noncalendar filings into consideration, the 1996 statistics represent all corporations filing returns with accounting periods ending during the period from July 1996 to June 1997. Also, many corporations, including some of the largest, request 6-month filing extensions. The combination of noncalendar year filing and filing extensions means that the last returns due to be received by IRS for the Tax Year 1996 (those with accounting periods ending in June 1997, which must therefore be filed by October 1997) could be timely filed as late as March 1998, if the 6-month extension of the October 1997 due date is taken into account. Normal administrative processing time lags required that the sampling process remain open for the 1996 study until June 30, 1998. However, a few very large returns for Tax Year 1996 were added to the sample as late as November 1998.

Each corporation is assigned a permanent and unique Employer Identification Number (EIN). The EIN is used as the basis for random selection. A pseudo-random number (PRN) is generated using the EIN as the seed. The last four digits of the PRN, called the transformed taxpayer identification number (TTIN), are compared to the sampling rates; a corporation for which the value of its TTIN is below the sampling rate multiplied by 10,000 is selected in the sample. The algorithm for generating the TTIN does not change from year to year. Consequently, any corporation selected into the sample in a given year will be selected again the next year, providing that the corporation files a return using the same EIN in the two years and that it falls into a stratum with the same or higher rate. If the corporation falls into a stratum with a lower rate, the chance of selection will correspond to the ratio of the second year to the first year selection probabilities. If the corporation files with a new EIN, the probability of being selected will be independent of the prior year selection. Due to the fact that corporations typically maintain the same EINs, this use of the EIN for the basis of sample selection results in many of the

same corporations selected into the sample from year to year. This also results in a reduction of the sample variance for estimates of year-to-year change [2].

Data Capture

Data processing for SOI begins with information already extracted for administrative purposes; over 100 items are available from the BMF system for nonconsolidated Form 1120 returns. Some 900 additional items are extracted from the tax returns during SOI processing. The administrative data are checked and corrected as necessary. The SOI data capture process can take as little time as fifteen minutes for a small, single entity corporation filing on Form 1120-A, or as long as a week for a large consolidated corporation filing several hundred attachments and schedules with the return. The process is further complicated by several factors:

- ‡ The 900 separate data items that may be extracted from any given tax return often require totals to be constructed from various other items on other parts of the return.
- ‡ Each 1120 form type has a different layout with different types of schedules and attachments, making data extraction less than uniform for the various form types.
- ‡ There is no legal requirement that a corporation meet its tax return filing requirements by filling in, line for line, the entire U.S. tax return form. Therefore, many corporate taxpayers report many of their financial details in schedules of their own design.
- ‡ There is no single accepted method of corporate accounting used throughout the country, but rather several accepted accounting "guidelines," many of which are unique to geographic locations. SOI attempts to standardize these differences during data abstraction and editing.
- ‡ Different companies may report the same data item, such as other current liabilities, on different lines of the tax form. Again, SOI attempts to standardize these differences.

In order to help overcome these complexities and differences due to taxpayer reporting, SOI prepares detailed instructions for the SOI editing unit at designated service centers each tax year. For Tax Year 1996, these instructions consisted of more than 800 pages covering normal and straightforward procedures and instructions for exceptions and nonstandard situations that might be encountered.

1996 Corporation Returns - Sample Description and Data Limitations

Data Cleaning

Statistical processing of the corporate returns took place in an online computer environment. This means that the data from returns were entered directly into the corporation database. In this context, the term "editing" refers to the combined interactive processes of data extraction, consistency testing, and error resolution. There are over 800 of these tests, which look for such inconsistencies as:

- | Impossible conditions, such as incorrect tax data for a particular form type;
- | Internal inconsistencies, such as items not adding to totals;
- | Questionable values, such as a bank with an unusually large amount reported for cost of goods sold and/or operations; and
- | Improper sample class codes, such as when a return has \$10,000 in total assets, but was selected as though it had \$1 million.

Data Completion

In addition to the tests mentioned above, missing data problems must be addressed and returns that are to be excluded from the tabulations must be identified. The data completion process focuses on these issues.

If the missing data items are from the balance sheet, then imputation procedures are used. If data for a whole return are missing because the return is unavailable to SOI during the data capture process, then, again, imputation procedures are used in certain cases.

A ratio-based imputation procedure is used to estimate missing balance sheet items for all 1120 forms except those with less than 12-month accounting periods. The ratios are determined by the corporation's 1995 return if it is available; otherwise, the 1994 aggregate data for the corporation's minor industrial group are used. If the reported items in the balance sheet do not balance (i.e., the sum of asset items does not equal the sum of liability and shareholders' equity items), then missing items are imputed. If the total assets amount is among the missing items, this item is imputed first based on the ratio of total assets to business receipts (or total receipts) from either the corporation's 1995 return, or the 1994 aggregate data for the corporation's minor industry. The other missing asset and liability items are then imputed based on the ratios so that the total of all asset items and the total of all liability items are both equal to the total assets amount, whether this amount was reported or imputed. A detailed description of the balance sheet imputation process is given in

reference [3]. The following table shows the number of sampled returns that had balance sheet items imputed for Tax Years 1993 through 1996.

Tax Year	1993	1994	1995	1996
No. of Returns	214	230*	131*	154*

* Starting in Tax Year 1994, 1504(c) returns are counted as one return rather than separate entities when computing the number of imputed balance sheets.

For Tax Year 1996, of the 154 returns, 25 of them have imputed total assets, and the imputed total asset amount constitutes approximately .0011 percent of the estimated total assets of the active corporations in 1996.

Data for unavailable critical corporations are imputed in various ways, depending on what information is available at the time the SOI database is produced. Critical corporations include corporations with total assets greater than or equal to 5 percent of the total assets for the minor industrial group in which they are classified, and corporations for which total assets are over a specified limit which is dependent on the form type or the major industry. For critical corporations selected for the sample but unavailable for statistical processing, taxpayer-surveyed data are used. There are two such returns in the Tax Year 1996 data. For the critical corporations not selected for the sample, if the current tax return is not found in any of the IRS service centers and no other current tax data are available, data from the previous year's return are used with adjustments for tax law changes. There are two prior year returns in the Tax Year 1996 data.

Another part of the data-cleaning process is identifying sampled returns that are not used in the tabulation. The BMF system, used for sample selection, can include duplicate tax returns and other out-of-scope returns, such as returns for nonprofit corporations and prior-year tax returns. These include the following types of returns:

- | Inactive corporation returns (having neither current income nor deductions);
- | Duplicate returns;
- | Amended returns not associated with the original returns;
- | Tentative returns not associated with the revised returns;
- | Corporations exempt under Code section 931;
- | Corporations exempt under Code section 1247;

1996 Corporation Returns - Sample Description and Data Limitations

- | Corporations exempt under Section 883 of the IRC;
- | "Cost corporation" returns exempt under Revenue Ruling 52-542;
- | Corporations exempt under Code section 501(c)(15);
- | Nonresident foreign corporations having no income effectively connected with a trade or business within the United States;
- | U.S. Virgin Island corporations exempt under Code section 934;
- | Political organizations filing under Code section 527;
- | General stock ownership corporations exempt from tax;
- | Homeowners' associations under Code section 528;
- | Information returns reporting no tax because of tax treaty or convention according to Code section 894;
- | Most prior-year returns with total assets under \$250 million filed on tax forms for years prior to 1995 and with accounting periods ending before July 1996;
- | Returns filed on a form type which should not be included in the SOI sample;
- | Fraudulent returns;
- | Returns of businesses incorporated in a tax-exempt U.S. Possession.

The following table displays the number of sampled returns that were excluded from tabulations and the percentages they represent of the total sample sizes in Tax Years 1993 through 1996.

Type of Return	Tax Year			
	1993	1994	1995	1996
Inactive	1,188	1,367	1,466	1,070
Duplicate	166	634	984	653
Other*	2,958	2,009	2,217	1,512
Total	4,312	4,010	4,667	3,235
% of Sample	4.71	4.22	4.78	3.44

* Includes prior-year returns.

Estimates of the number of active corporations by form type for Tax Years 1993 through 1996 are provided in the next table.

Form Type	Tax Year			
	1993	1994	1995	1996
1120	1,775,931	2,038,870	2,043,818	2,062,341
1120-A	265,627	257,125	257,439	241,536
1120S	1,901,505	2,023,754	2,153,119	2,304,416
1120-L	1,876	1,775	1,646	1,725
1120-PC	2,623	2,674	2,789	3,435
1120-RIC	6,796	7,519	8,201	8,541
1120-REIT	346	393	465	526
1120-F	9,925	10,259	6,690*	8,849
Total	3,964,629	4,342,368	4,474,167	4,631,370

Note: Detail may not add to total due to rounding.

* This estimate is significantly lower than in previous years (see section on Coverage Errors).

Estimation

The estimates produced in this report of the total number of corporations and associated money amounts are based on weighted sample results. Either a one-step process or a two-step process was used to determine the weights, depending on the return's form type.

Under the one-step process, the weights are assigned as the reciprocal of the achieved sample rate. These weights are used to produce the aggregated total frequencies and money amounts published in this report for Forms 1120-F, 1120-L, 1120-PC, 1120-RIC, 1120-REIT and Form 1120 with Form 5735 attached.

The two-step process was used to improve the industry estimates. The first stage is identical to the one-step process as described above and provides an initial weight for the return. The second stage involves poststratification by industry. During poststratification, certain cells have small sample sizes. To handle this problem, a raking ratio estimation approach is applied during poststratification in order to determine the final weights [4]. Restrictions are placed on the raking process to produce final weights that fall within the range $1/(2/3) \times$ original weight to $1/(3/2) \times$ original weight. These final weights are used to produce the aggregated frequencies and money amounts published in this report for Forms 1120, 1120-A and 1120S.

Data Limitations and Measures of Variability

Several extensive quality review processes were used to improve the quality of the data. The review processes began at the sample selection stage with weekly monitoring of the sample to ensure that the proper number of returns was being selected. They continued through the data collection, data cleaning,

1996 Corporation Returns - Sample Description and Data Limitations

and data completion procedures with consistency testing. Part of the review process included extensive comparisons between the 1996 data and the 1995 data. A great amount of effort was made at every stage of processing to ensure data integrity.

Sampling Error

Since the corporation estimates are based on a sample, they may differ from figures that would have been obtained if a complete census of all income tax returns had been taken. The particular sample used to produce the results in this report is one of a large number of possible samples that could have been selected under the same sample design. Estimates derived from one of the possible samples could differ from those derived for any other sample, and from the population aggregates. The deviation of a sample estimate from the average of all possible similarly selected samples is called the sampling error. The standard error (SE) is a measure of the average magnitude of the sampling errors over all possible samples.

The standard error is the most commonly used measure of the sampling error and can be estimated from the sample. Sometimes, for convenience, the standard error is expressed as a percentage of the

value being estimated. This is called the coefficient of variation (CV) of the estimate. The coefficient of variation can be used in assessing the reliability of an estimate.

The coefficient of variation of an estimate is calculated by dividing the standard error by the estimate. Coefficients of variation by industrial groupings for the estimated number of returns, as well as for selected money amount estimates, are shown in Table 1 beginning on page 29. For the estimated number of returns by asset size and industrial division, coefficients of variation are given in Figure D.

The coefficient of variation, CV(X), can be used to construct confidence intervals of the estimate X. The standard error, which is required for the confidence interval, must first be calculated. For example, the estimated number of manufacturing companies with net income and its coefficient of variation can be found in Table 1 and used to calculate the standard error:

$$\begin{aligned} SE(X) &= X \cdot CV(X) \\ &= 191,254 \times 2.54/100 \\ &= 4,858 \end{aligned}$$

Figure D--CVs for Number of Returns, by Asset Size and Industrial Division, Tax Year 1996

Industrial division	All Asset Sizes	Size of total assets				
		Zero Assets	\$1 Under \$100,000	\$100,000 under \$250,000	\$250,000 under \$500,000	\$500,000 under \$1,000,000
	(1)	(2)	(3)	(4)	(5)	(6)
All Industries ¹	0.20	3.19	0.53	0.80	0.73	0.53
Agriculture	2.62	21.00	5.81	5.17	3.55	2.83
Mining	7.31	29.43	13.51	20.00	13.59	10.50
Construction	1.17	10.46	2.30	3.27	2.91	2.11
Manufacturing	1.87	13.33	4.51	5.30	3.49	2.57
Transportation	2.60	15.56	4.52	5.81	5.31	5.23
Wholesale and retail trade	0.68	7.03	1.53	1.58	1.46	1.19
Finance, insurance, and real estate	0.89	7.15	2.06	2.37	1.88	1.52
Services	0.66	5.47	1.04	1.97	2.26	2.18
Industrial division	Size of total assets—continued					
	\$1,000,000 under \$5,000,000	\$5,000,000 under \$10,000,000	\$10,000,000 Under \$25,000,000	\$25,000,000 under \$50,000,000	\$50,000,000 under \$100,000,000	\$100,000,000 under \$250,000,000
	(7)	(8)	(9)	(10)	(11)	(12)
All Industries ¹	0.27	0.59	0.45	0.57	0.04	0.03
Agriculture	1.96	5.34	4.97	6.54	1.11	0.97
Mining	5.14	6.40	5.14	5.38	0.84	0.66
Construction	1.07	2.27	2.02	3.13	0.59	0.70
Manufacturing	1.17	1.80	1.13	1.21	0.20	0.17
Transportation	2.52	4.84	2.79	3.41	0.50	0.40
Wholesale and retail trade	0.64	1.17	0.89	1.31	0.26	0.23
Finance, insurance, and real estate	0.88	1.77	1.18	1.19	0.10	0.06
Services	1.37	2.93	2.01	2.29	0.34	0.31

¹Includes returns not allocable by industrial division.

1996 Corporation Returns - Sample Description and Data Limitations

Assume that a 95-percent confidence interval for the number of returns in manufacturing is desired. The 95- percent confidence interval is constructed as follows:

$$\begin{aligned} X \pm 2SE(X) &= 191,254 \pm (2 \times 4,858) \\ &= 191,254 \pm 9,716 \end{aligned}$$

Thus, the interval estimate is 181,538 returns to 200,970 returns. This means that if all possible samples were selected under essentially the same general conditions and using the same sample design, and if an estimate and its standard error were calculated from each sample, then approximately 95 percent of the intervals from two standard errors below the estimate to two standard errors above the estimate would include the average estimate derived from all possible samples. Thus, for a particular sample, it can be said with 95-percent confidence that the average of all possible samples is included in the constructed interval. This average of the estimates derived from all possible samples would be equal to or near the value obtained from a census.

Nonsampling Error

In addition to sampling error, nonsampling error can also affect the estimates. Nonsampling errors can be classified into two groups: random errors whose effects may cancel out and systematic errors whose effects tend to remain somewhat fixed and result in bias.

Nonsampling errors can be categorized as coverage errors, nonresponse errors, processing errors, or response errors. These errors can be the result of the inability to obtain information about all returns in the sample, differing interpretations of tax concepts or instructions by the taxpayer, inability of a corporation to provide accurate information at the time of filing (data are collected before auditing), inability to obtain all tax schedules and attachments, errors in recording or coding the data, errors in collecting or cleaning the data, errors made in estimating for missing data, and failure to represent all population units.

Coverage Errors

Coverage errors in the SOI Corporation data can result from the difference between the time frame for sampling and the actual time needed for filing and processing the returns. As stated above, many of the largest corporations receive extensions to their filing periods and, as a result, may file their returns after sample selection has ended for that tax year. However, any of the largest returns found are added into the file until the final file is produced.

Coverage problems within industrial divisions in the SOI Corporation study result from the way

consolidated returns may be filed. The Internal Revenue Code permits a parent corporation to file a single return, which includes the combined financial data of the parent and all its subsidiaries. These data are not separated into the different industries but are entered only into the industry with the largest receipts. Thus, there is undercoverage of financial data within certain industries and overcoverage in others. Coverage problems within industrial divisions present a limitation on any analysis done with the sample results.

In 1996, as in 1994 and 1995, there was a processing problem prior to the sampling operation which resulted in some Form 1120-F returns filed by corporations with income "effectively connected with a U.S. trade or business" being excluded from the sampling frame. Specifically, these returns were incorrectly coded as not having effectively connected income. This resulted in undercoverage of the Form 1120-F population.

To overcome the undercoverage problem beginning with Tax Year 1997, all Form 1120-F returns regardless of their coding will be subject to sampling, thus ensuring that all those with effectively connected income are definitely included in the sampling frame. A preliminary estimate from the 1997 file does show a marked increase over the 1994, 1995, and 1996 Form 1120-F active population estimates; this 1997 population estimate is much more consistent with the estimates from SOI years prior to 1994.

Presently, SOI is engaged in researching various statistical methodologies to try to improve the estimates (number of returns and money amounts) for the affected years. Currently, two avenues are being pursued. The first approach being examined is to perform statistical adjustments using either a logistic regression or a time-series analysis. Another approach being pursued is obtaining returns that were initially omitted from the sample due to the undercoverage. We have identified a sample from the 1995 Form 1120-F population of returns coded as not having effectively connected income similarly to what is described above for the 1997 study. We will identify the returns that have effectively connected income and use this additional sample, along with the sample previously selected, in order to compute adjusted weights. These derived estimates should be close to what would have been obtained had the returns not been omitted, since we are in effect re-creating the complete 1995 sample.

Nonresponse Errors

Unit nonresponse for SOI occurs when a sampled return is unavailable for SOI processing. For example, other areas of the IRS such as

1996 Corporation Returns - Sample Description and Data Limitations

Examination, Collection, or District Offices may have the return at the time the return is needed for statistical processing. These returns are termed "unavailable returns." In 1996, there were 104 unavailable returns in the corporation study, which constituted about .11 percent of the total sample size. The following table gives the number of unavailable returns and their percentages of total sample sizes for Tax Years 1993 through 1996.

Tax Year	1993	1994	1995	1996
No. of Returns	118	113	138	104
% of Sample	0.13	0.12	0.14	0.11

Processing Errors

Errors in recording, coding or processing the data can cause a return to be sampled in the wrong sampling class. This type of error is called a mis-stratification error. One example of how a return might be mis-stratified is the following: a corporation files a return with total assets of \$10,000.23 and net income of \$5,000.00. A processing error causes the cents to be keyed in as dollars so that the return is classified according to total assets of \$1,000,023 and net income of \$5,000. The return would be mis-stratified according to the incorrect value of total assets.

The following table shows the number of mis-stratified returns for Tax Years 1993 through 1996.

Tax Year	1993	1994	1995	1996
No. of Returns	1,082	1,324	1,420	1,618

Mis-stratified returns in the sample were reclassified into their proper sampling classes after complete data capture. The population of returns that needed to be reclassified was estimated from the sample and the stratum population sizes were adjusted accordingly [5]. Population and sample totals were minimally affected by reclassification, and an analysis of the sample results tended to confirm that mis-stratified returns occurred randomly. Steps are being taken by both the Centers and the SOI Division to minimize the number of mis-stratified returns.

Response errors

Response errors are due to data being captured before auditing. Some purely arithmetical errors

made by the taxpayer are corrected during the data capture and cleaning processes. Because of time constraints, adjustments to a return during auditing are not incorporated into the SOI file.

Industrial Classification

The industry classification used in this report generally conforms to the former Enterprise Standard Industrial Classification (ESIC) authorized by The Office of Information and Regulatory Affairs in The Office of Management and Budget (OMB). This classification was designed to classify companies, which are often engaged in more than one industry activity, into only one industry category. The structure of this classification follows closely along the lines of the underlying Standard Industrial Classification (SIC) Manual, also authorized by OMB, which is designed as a means of classifying establishments. Some departures from the ESIC system were made by SOI for financial industries in order to reflect particular provisions of the Internal Revenue Code.

References

- [1] Jones, H. W., and McMahon, P. B. (1984), "Sampling Corporation Income Tax Returns for Statistics of Income, 1951 to Present," *1984 Proceedings of the Section on Survey Research Methods*, American Statistical Association, pp. 437-442.
 - [2] Harte, J. M. (1986), "Some Mathematical and Statistical Aspects of the Transformed Taxpayer Identification Number: A Sample Selection Tool Used at IRS," *1986 Proceedings of the Section on Survey Research Methods*, American Statistical Association, pp. 603-608.
 - [3] Überall, B. (1995), "Imputation of Balance Sheets for the 1992 SOI Corporate Program," *1995 Proceedings of the Section on Survey Research Methods*, American Statistical Association, pp. 275-280.
 - [4] Oh, H. L. and Scheuren, F. J. (1987), "Modified Raking Ratio Estimation," *Survey Methodology*, Statistics Canada, Vol. 13, No. 2, pp. 209-219.
 - [5] Mulrow, J. M. and Jones, H. W. (1989), "Sampling Administrative Records: Detection and Correction of Stratification Errors," *Statistics of Income and Related Administrative Record Research: 1988-1989*, Internal Revenue, December 1990, pp. 139-144.
- SOURCE: IRS, Statistics of Income 1996 Corporation Income Tax Returns, Publication 16 (rev. 10-99)