

Section 2

Description of the Sample

This section describes the sample design and selection, the method of estimation, the sampling variability of the estimates, and the methodology of computing confidence intervals.

Domain of Study

The statistics in this report are estimates from a probability sample of unaudited Individual Income Tax Returns, Forms 1040, 1040A, and 1040EZ (including electronic returns) filed by U.S. citizens and residents during Calendar Year 2006.

All returns processed during 2006 were subjected to sampling except tentative and amended returns. Tentative returns were not subjected to sampling because the revised returns may have been sampled later, while amended returns were excluded because the original returns had already been subjected to sampling. A small percentage of returns were not identified as tentative or amended until after sampling. These returns, along with those that contained no income information, were excluded in calculating estimates. This resulted in a small difference between the population total (134,494,440 returns) reported in Table C and the estimated total of all returns (134,372,678)

reported in other tables.

The estimates in this report are intended to represent all returns filed for Tax Year 2005. While most of the returns processed during Calendar Year 2006 were for Tax Year 2005, the remaining returns were mostly for prior years, and a few for non-calendar years ending during 2006 and 2007. Returns for prior years were used in place of 2005 returns received and processed after December 31, 2006. This was done based on the assumption that the characteristics of returns due, but not yet processed, can best be represented by the returns for previous income years that were processed in 2006.

Sample Design and Selection

The sample design is a stratified probability sample, in which the population of tax returns is classified into subpopulations, called strata, and a sample is randomly selected independently from each stratum. Strata are defined by:

1. Nontaxable (including no alternative minimum tax) with adjusted gross income or expanded income of \$200,000 or more.

Valerie Testa, and Jana Scali designed the sample and prepared the text and tables in this section under the direction of Yahia Ahmed, Chief, Mathematical Statistics Section, Statistical Computing Branch.

2. High combined business and farm total receipts of \$50,000,000 or more.
3. Presence or absence of special Forms or Schedules (Form 2555, Form 1116, Form 1040 Schedule C, and Form 1040 Schedule F).
4. Indexed positive or negative income. Sixty variables are used to derive positive and negative incomes. These positive and negative income classes are deflated using the Chain-Type Price Index for the Gross Domestic Product to represent a base year of 1991. (See footnote 1 for details.)
5. Potential usefulness of the return for tax policy modeling. Thirty-two variables are used to determine how useful the return is for tax modeling purposes.

Table C shows the population and sample count for each stratum after collapsing some strata with the same sampling rates. (See references 1 and 2 for details.) The sampling rates range from 0.10 percent to 100 percent.

Tax data processed to the IRS Individual Master File at the Enterprise Computing Center at Martinsburg during Calendar Year 2006 were used to assign each taxpayer's record to the appropriate stratum and to determine whether or not the record should be included in the sample. Records are selected for the sample either if they possess certain combinations of the four ending digits of the social security number, or if their ending five digits of an eleven-digit number generated by a mathematical transformation of the SSN is less than or equal to the stratum sampling rate times 100,000. (See reference 3 for details.)

Data Capture and Cleaning

Data capture for the SOI sample begins with the designation of a sample of administrative records. While the sample was being selected, the process was continually monitored for sample selection and data collection errors. In addition, a small subsample of returns was selected and independently reviewed, analyzed, and processed

for a quality evaluation.

The administrative data and controlling information for each record designated for this sample was loaded onto an online database at the Cincinnati Submission Processing Center. Computer data for the selected administrative records were then used to identify inconsistencies, questionable values, and missing values as well as any additional variables that an editor needed to extract for each record. The editors use a hardcopy of the taxpayer's return to enter the required information onto the online system.

After the completion of service center review, data were further validated, tested, and balanced. Adjustments and imputations for selected fields based on prior year data and other available information were used to make each record internally consistent. Finally, prior to publication, all statistics and tables were reviewed for accuracy and reasonableness in light of provisions of the tax law, taxpayer reporting variations and limitations, economic conditions, and comparability with other statistical series.

Some returns designated for the sample were not available for SOI processing because other areas of IRS needed the return at the same time. For Tax Year 2005, 0.10 percent of the sample returns were unavailable.

Method of Estimation

Weights were obtained by dividing the population count of returns in a stratum by the number of sample returns for that stratum. The weights were adjusted to correct for misclassified returns. These weights were applied to the sample data to produce all of the estimates in this report.

Sampling Variability and Confidence Intervals

The sample used in this study is one of a large number of samples that could have been selected using the same sample design. The estimates calculated from these different samples would vary. The standard error (SE) of an estimate is a measure of the variation among the estimates from the possible samples and, thus, is a measure of the

precision with which an estimate from a particular sample approximates the average of the estimates calculated from all possible samples.

The standard error may be expressed as a percentage of the value being estimated. This ratio is called the coefficient of variation (CV). Tables 1.4 CV, 2.1 CV, and 3.3 CV contain estimated CV's for the estimates included in Tables 1.4, 2.1, and 3.3 of this report.

The sample estimate and an estimate of its standard error permit the construction of interval estimates with prescribed confidence that the interval includes the population value. If all possible samples were selected under essentially the same conditions and an estimate and its estimated standard error were calculated from each sample, then:

1. About 68 percent of the intervals from one standard error below the estimate to one standard error above the estimate would include the population value. This is a 68 percent confidence interval.
2. About 95 percent of the intervals from two standard errors below the estimate to two standard errors above the estimate would include the population value. This is a 95 percent confidence interval.

For example, from Table 1.4, the estimate for State Income Tax Refunds, X, is \$22.205 billion, and its related coefficient of variation, CV(X), is 0.77 percent. The standard error of the estimate, SE(X), needed to construct the confidence interval estimate, is:

$$\begin{aligned} SE(X) &= X \cdot CV(X) \\ &= (\$22.205 \times 10^9) \cdot (0.0077) \\ &= \$0.171 \text{ billion} \end{aligned}$$

The p percent confidence interval is calculated using the formula:

$$X \pm z \cdot SE(X)$$

where z takes the value 1, 2, or 3 when p is 68, 95,

or 99, respectively. Based on these data, the 68 percent confidence interval is from \$22.034 billion to \$22.376 billion, the 95 percent confidence interval is from \$21.863 billion to \$22.547 billion, and the 99 percent confidence interval is from \$21.692 billion to \$22.718 billion.

Table Presentation

Whenever a weighted frequency is less than 3, the estimate and its corresponding amount are combined or deleted in order to avoid disclosure of information for specific taxpayers. (The combined or deleted data, if any, are included in the corresponding column totals.) These combinations and deletions are indicated by a double asterisk (**). Estimates based on less than 10 sampled returns are considered to be unreliable. These estimates are noted by a single asterisk (*) to the left of the data unless all of the sampled returns are selected with certainty (at the 100 percent rate).

In the tables, a dash (-) in place of a frequency or an amount indicates that either no returns in the population had the characteristic or the characteristic was so rare that it did not appear on any of the sampled returns.

Footnote

[1] Indexing of positive and negative income is done by dividing each by the ratio of the Chain-Type Price Index for the Gross Domestic Product for the fourth quarter of 2004 to the fourth quarter of the base year of 1991. The indices were calculated using the Gross Domestic Product (GDP) Chain-type Price Index found in the table titles "Quantity and Price Indexes for Gross Domestic Product" released to the public on November 30, 2005 on the BEA web site (<http://www.bea.doc.gov/>).

References

[1] Hostetter, S., Czajka, J. L., Schirm, A. L., and O'Connor, K. (1990), "Choosing the Appropriate Income Classifier for Economic Tax Modeling," in *Proceedings of the Section*

on *Survey Research Methods*, American Statistical Association, 419-424.

[2] Schirm, A. L., and Czajka, J. L. (1991), "Alternative Designs for a Cross-Sectional Sample of Individual Tax Returns: the Old and the New," *Proceedings of the Section on Survey Research Methods*, American Statistical Association, 163-168.

[3] Harte, J.M. (1986), "Some Mathematical and Statistical Aspects of the Transformed Taxpayer Identification Number: A Sample Selection Tool Used at IRS," *Proceedings of the Section on Survey Research Methods*, American Statistical Association, 603-608.

Table C.—Number of Individual Income Tax Returns in the Population and Sample by Sampling Strata for 2005

Description of the sample strata	Description of the sample strata												Number of returns	
	Degree of interest ²	Form 1040, with Form 1116 or Form 2555		Form 1040, with Schedule C but without Form 1116 or Form 2555		Form 1040, with Schedule F but without Schedule C, Form 1116 or Form 2555		Form 1040, with other Schedules and Forms and Forms 1040A and 1040EZ		Population counts ¹	Sample counts			
		(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)		
Grand total		4,284,788	61,013	20,813,190	57,385	1,407,010	5,671	107,976,672	156,137	134,494,440	292,966			
Form 1040 returns only with adjusted gross income or expanded income of \$200,000 and over, with no income tax after credits and no additional tax for tax preferences, total										12,550	12,550			
Form 1040 returns only with combined Schedule C (business or profession) total receipts of \$50,000,000 and over, total										230	230			
Other Returns, total										134,481,660	280,186			
Number of Returns by type of form attached														
Description of the sample strata	Degree of interest ²	Form 1040, with Form 1116 or Form 2555		Form 1040, with Schedule C but without Form 1116 or Form 2555		Form 1040, with Schedule F but without Schedule C, Form 1116 or Form 2555		Form 1040, with other Schedules and Forms and Forms 1040A and 1040EZ		Population counts	Sample counts			
		(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)					
Total		4,284,788	61,013	20,813,190	57,385	1,407,010	5,671	107,976,672	156,137					
Indexed Negative Income ³														
Under \$30,000	1	197,051	204	2,611,879	2,650	89,378	84	30,444,833	30,497	30,444,833	30,497			
\$30,000 under \$10,000,000	2	230,222	361	4,408,123	6,695	133,661	211	26,046,587	25,913	28,944,895	28,851			
\$10,000,000 under \$5,000,000	3-4	402,885	419	1,925,464	1,893	166,353	178	5,460,294	8,393	10,232,300	15,660			
\$5,000,000 under \$2,000,000	1-2	452,589	690	3,660,439	5,617	243,641	390	21,248,325	21,321	23,743,027	23,811			
\$2,000,000 under \$60,000	3-4	691,802	693	2,200,594	2,249	220,474	223	5,898,461	9,454	10,255,130	16,151			
\$60,000 under \$120,000	1-3	532,957	816	2,623,422	4,003	190,337	263	10,729,825	10,609	13,842,695	13,774			
\$120,000 under \$250,000	4	218,697	429	365,057	737	86,499	170	2,999,844	4,478	6,346,560	9,560			
\$250,000 under \$500,000	1-3	643,231	2,178	1,449,248	4,782	90,347	268	1,073,621	2,151	1,744,074	3,487			
\$500,000 under \$1,000,000	4	412,531	2,943	519,045	3,819	67,059	460	1,904,356	6,345	4,087,182	13,573			
\$1,000,000 under \$2,000,000	All	199,495	4,871	152,260	3,934	20,409	505	626,369	4,433	1,625,004	11,655			
\$2,000,000 under \$5,000,000	All	85,879	10,383	40,305	4,970	5,384	614	177,374	4,307	549,538	13,617			
\$5,000,000 under \$10,000,000	All	43,239	14,071	13,401	4,405	1,720	559	52,923	6,514	184,491	22,481			
\$10,000,000 under \$50,000,000	All	11,924	2,676	2,676	4,058	340	340	19,351	6,189	77,711	25,224			
\$50,000,000 under \$10,000,000	All	8,285	8,285	1,235	1,235	156	156	4,058	4,058	18,998	18,998			
\$10,000,000 or more	All							1,959	1,959	11,635	11,635			

¹ This population includes an estimated 963,940 returns that were excluded from other tables in this report because they contained no income information or tentative returns identified after sampling.
² Each population member is assigned a degree of interest based on how useful it is for tax modeling purposes. Degree of interest ranges from one (1) to four (4), with a one being assigned to returns that are the least interesting, and a four being assigned to those that are the most interesting. 'All' refers to income classes for which returns with all four degrees of interest are assigned.
³ Positive and Negative Income classes are divided by a Chain-Type Price Index for the Cross Domestic Product of 1.2510 to represent a base year of 1991.