

## 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 2011.

All returns processed during 2011 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

(143,170,763 returns) reported in Table C and the estimated total of all returns (142,892,051) reported in other tables.

The estimates in this report are intended to represent all returns filed for Tax Year 2010. While most of the returns processed during Calendar Year 2011 were for Tax Year 2010, the remaining returns were mostly for prior years, and a few for non-calendar years ending during 2009 and 2010. Returns for prior years were used in place of 2010 returns received and processed after December 31, 2011. 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 2011.

### 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 designed the sample and prepared the text and tables in this section under the direction of Tammy Rib, Chief, Mathematical Statistics Section, Statistical Computing Branch.*

2. High business 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 2011 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 2010, 0.03 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 \$27.454 billion, and its related coefficient of variation, CV(X), is 0.76 percent. The standard error of the estimate, SE(X), needed to construct the confidence interval estimate, is:

$$\begin{aligned} SE(X) &= X \cdot CV(X) \\ &= (\$27.454 \times 10^9) \cdot (0.0076) \\ &= \$0.209 \text{ billion} \end{aligned}$$

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 \$27.245 billion to \$27.663 billion, the 95 percent confidence interval is from \$27.036 billion to \$27.872 billion, and the 99 percent confidence interval is from \$26.827 billion to \$28.081 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 2009 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 "Price Indexes for Gross Domestic Product" released to the public on November 23, 2010 on the BEA web site (<http://www.bea.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 2010**

Description of the sample strata	Degree of interest [2]	Description of the sample strata												Number of returns			
		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
		(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)	(15)	(16)	
<b>Grand total</b>		5,486,531	60,996	22,187,594	52,353	1,353,840	6,050	114,108,428	155,177	143,170,763	308,946	34,072	298	143,136,393	274,576		
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																	
Form 1040 returns only with combined Schedule C (business or profession) total receipts of \$50,000,000 and over, total																	
Other Returns, total																	
		Number of Returns by type of form attached															
		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
		(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)	(13)	(14)	(15)	(16)	(17)
Total		5,486,531	60,996	22,187,594	52,353	1,353,840	6,050	114,108,428	155,177	143,170,763	308,946	34,072	298	143,136,393	274,576		
Indexed Negative Income [3]																	
Under \$10,000 or more	All	517	517	1,062	1,062	154	154	129	129	1,291	1,291	3,024	3,024	3,024	3,024		
\$5,000,000 under \$10,000,000	All	983	983	1,780	1,780	262	262	2,300	2,300	5,325	5,325	5,325	5,325	5,325	5,325		
\$2,000,000 under \$5,000,000	All	4,305	1,450	6,833	2,326	1,061	384	8,895	3,050	21,094	7,210	7,210	7,210	7,210	7,210		
\$1,000,000 under \$2,000,000	All	9,143	1,468	13,962	2,258	2,465	380	18,024	2,796	43,594	6,902	6,902	6,902	6,902	6,902		
\$500,000 under \$1,000,000	All	21,100	692	33,521	1,150	6,177	1,309	42,072	1,309	102,870	3,368	3,368	3,368	3,368	3,368		
\$250,000 under \$500,000	All	43,545	442	74,355	724	12,242	127	95,691	928	225,833	2,221	2,221	2,221	2,221	2,221		
\$120,000 under \$250,000	All	81,478	389	148,710	750	20,028	106	208,997	1,016	459,213	2,261	2,261	2,261	2,261	2,261		
\$60,000 under \$120,000	All	93,530	245	184,361	605	21,472	78	303,026	903	602,389	1,831	1,831	1,831	1,831	1,831		
Under \$60,000	All	71,405	130	399,823	733	30,464	54	899,766	1,698	1,401,458	2,615	2,615	2,615	2,615	2,615		
Indexed Positive Income [3]																	
Under \$30,000	1																
Under \$30,000	2	237,923	252	3,554,726	3,491	83,557	90	30,017,317	30,092	33,893,523	33,925	33,925	33,925	33,925	33,925		
Under \$30,000	3-4	207,308	189	5,293,117	5,447	107,329	117	6,925,044	6,802	12,532,798	12,555	12,555	12,555	12,555	12,555		
\$30,000 under \$60,000	1-2	610,028	593	1,860,689	1,760	161,107	170	21,171,494	21,229	23,803,318	23,752	23,752	23,752	23,752	23,752		
\$30,000 under \$60,000	3-4	539,597	550	3,783,108	3,756	247,918	224	6,685,816	6,849	11,256,439	11,379	11,379	11,379	11,379	11,379		
\$60,000 under \$120,000	1-3	956,344	963	2,110,530	2,190	204,511	227	10,783,313	10,659	14,054,698	14,039	14,039	14,039	14,039	14,039		
\$60,000 under \$120,000	4	674,901	677	2,457,716	2,457	181,812	158	3,072,233	3,072	6,386,662	6,364	6,364	6,364	6,364	6,364		
\$120,000 under \$250,000	1-3	279,031	944	338,505	1,069	76,031	282	1,089,010	3,665	1,782,577	5,980	5,980	5,980	5,980	5,980		
\$120,000 under \$250,000	4	811,225	2,625	1,312,549	4,453	91,318	280	2,029,537	6,765	4,244,629	14,123	14,123	14,123	14,123	14,123		
\$250,000 under \$500,000	All	487,617	3,589	442,429	3,170	70,984	505	601,204	4,288	1,602,234	11,552	11,552	11,552	11,552	11,552		
\$500,000 under \$1,000,000	All	217,104	5,379	124,815	3,141	26,641	645	150,828	3,736	519,388	12,901	12,901	12,901	12,901	12,901		
\$1,000,000 under \$2,000,000	All	84,721	10,217	32,337	4,050	6,484	788	39,410	4,822	162,952	19,877	19,877	19,877	19,877	19,877		
\$2,000,000 under \$5,000,000	All	38,461	12,437	9,919	3,214	1,502	481	12,967	4,254	62,849	20,386	20,386	20,386	20,386	20,386		
\$5,000,000 under \$10,000,000	All	9,915	9,915	1,920	1,920	229	229	2,473	2,473	14,537	14,537	14,537	14,537	14,537	14,537		
\$10,000,000 or more	All	6,350	6,350	827	827	92	92	1,264	1,264	8,533	8,533	8,533	8,533	8,533	8,533	8,533	

[1] This population includes an estimated 278,712 returns that were excluded from other tables in this report because they contained no income information or represented amended 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 Gross Domestic Product of 1.4530 to represent a base year of 1991.

