

# Accuracy of the Data (2005)

## INTRODUCTION

The data contained in these data products are based on the American Community Survey (ACS) and Puerto Rico Community Survey (PRCS) sample interviewed from January 1, 2005 through December 31, 2005. [Unless otherwise specified, the term “ACS” in this document will refer to both the ACS and PRCS.] Beginning in 2005, the ACS sample expanded to include all counties and county-equivalents in the United States, and all municipios in Puerto Rico (PR). The ACS, like any other statistical activity, is subject to error. The purpose of this documentation is to provide data users with a basic understanding of the ACS sample design, estimation methodology, and accuracy of the ACS data. The ACS is sponsored by the U.S. Census Bureau, and is an integral part of the plan for the 2010 Census.

Additional information on the operational aspects of the ACS, including data collection and processing, can be found in the ACS Operations Plan (<http://www.census.gov/acs/www/Downloads/OpsPlanfinal.pdf>).

## DATA COLLECTION

The ACS and PRCS employ three modes of data collection:

- Mailout/Mailback
- Computer Assisted Telephone Interview (CATI)
- Computer Assisted Personal Interview (CAPI)

With the exception of addresses in Remote Alaska, the general timing of data collection is:

- Month 1: Addresses determined to be mailable are sent a questionnaire via the U.S. Postal Service.
- Month 2: All mail non-responding addresses with an available phone number are sent to CATI.
- Month 3: A sample of mail non-responses without a phone number, CATI non-responses, and unmailable addresses are selected and sent to CAPI.

All Remote Alaska addresses are assigned to one of two data collection periods, January-April, or September-December. Data for these addresses are collected using CAPI only. Note that mail responses are accepted during all three months of data collection.

## SAMPLING FRAME

The sampling frame for the ACS is created from the Master Address File (MAF), which is a database maintained by the Census Bureau containing a listing of residential and commercial addresses in the U.S. and Puerto Rico (PR). The MAF is updated twice each year with the Delivery Sequence Files provided by the U.S. Postal Service which cover only the U.S. These files identify mail drop points and provide the best available source of changes and updates to the housing unit inventory. The MAF is also updated with the results from various Census Bureau field operations, including the ACS.

## SAMPLE DESIGN

The ACS employs a two-stage, two-phase sample design. The ACS first-stage sample consists of two separate samples, Main and Supplemental, each chosen at different points in time. Together, these constitute the first-stage sample. Both the Main and the Supplemental samples are chosen in two phases referred to as first- and second-phase sampling. Subsequent to second-phase sampling, sample addresses are randomly assigned to one of the twelve months of the sample year. The second-stage of sampling occurs when the CAPI sample is selected (see Section 2 below).

The Main sample is selected during the summer proceeding the sample year. Approximately 99% of the sample is selected at this time. Each address in sample is randomly assigned to one of the 12 months of the sample year. Supplemental sampling occurs in January/February of the sample year and accounts for approximately 1% of the overall first-stage sample. The Supplemental sample is allocated to the last nine months of the sample year.

Several of the steps used to select the first-stage sample are common to both Main and Supplemental sampling. The descriptions of the steps included in the first-stage sample selection below indicate which are common to both and which are unique to either Main or Supplemental sampling.

### 1. First-Stage Sample Selection

- First-phase sampling (*performed during both Main and Supplemental sampling*) – First stage sampling defines the universe for the second stage of sampling through two steps. First, all addresses that were in a first-stage sample within the past four years are excluded from eligibility. This ensures that no address is in sample more than once in any five-year period. The second step is to select a 20% systematic sample of “new” units, i.e. those units that have never appeared on a previous MAF extract. Each new address is systematically assigned to either the current year or to one of four back-samples. This procedure maintains five equal partitions of the universe.
- Assignment of blocks to a second-phase sampling stratum (*performed during Main sampling only*) – Second-phase sampling uses seven distinct sampling rates in the U.S.

and five in PR. These rates are applied at a block level to addresses in the U.S. and PR by calculating a measure of size for each of the following entities:

- o Counties
- o Places (active, functioning governmental units)
- o School Districts (elementary, secondary, and unified)
- o American Indian Areas
- o Alaska Native Village Statistical Areas
- o Hawaiian Homelands
- o Minor Civil Divisions (MCDs) – in Connecticut, Maine, Massachusetts, Michigan, Minnesota, New Hampshire, New Jersey, New York, Pennsylvania, Rhode Island, Vermont, and Wisconsin (these are the states where MCDs are active, functioning governmental units)
- o Census Designated Places – in Hawaii only

The measure of size for all areas except American Indian and Alaska Native Village Statistical Areas is an estimate of the number of occupied housing units (HUs) in the area. This is calculated by multiplying the number of ACS addresses by the occupancy rate from Census 2000 at the block level. A measure of size for each Census Tract is also calculated in the same manner.

For American Indian and Alaska Native Village Statistical Areas the measure of size is the estimated number of occupied HUs multiplied by the proportion of people reporting American Indian or Alaska Native (alone or in combination) in Census 2000.

Each block is then assigned the smallest measure of size from the set of all entities it is a part of. These two measures are used to assign the first-stage sampling rates as shown in Table 1 below.

- Calculation of the second-phase sampling rates (*performed during Main sampling only*) – The sampling rates given in Table 1 are calculated using the distribution of ACS valid addresses by second-phase sampling stratum in such a way as to yield an overall target sample size for the year of approximately 3,000,000 in the U.S. and 36,000 in PR. These rates also account for expected growth of the HU inventory between Main and Supplemental of roughly 1%.
- Second-phase sample selection (*performed in Main and Supplemental*) – After each block is assigned to a second-phase sampling stratum, a systematic sample of addresses is selected from the second-phase universe (first-phase sample) within each county, county equivalent, and municipio.
- Sample Month Assignment (*performed in Main and Supplemental*) – After each phase of sampling, all sample addresses are randomly assigned to a sample month. Addresses selected during Main sampling are allocated to each of the 12 months. Addresses selected during Supplemental sampling are assigned to the months of April-December.

**Table 1. First-Stage Sampling Rate Categories for the United States and Puerto Rico**

Sampling Rate Category	Sampling Rates	
	United States	Puerto Rico
Blocks in smallest governmental units (MOS <sup>1</sup> < 200)	10.0%	10.0%
Blocks in smaller governmental units (200 ≤ MOS < 800)	6.9%	8.1%
Blocks in small governmental units (800 ≤ MOS ≤ 200)	3.6%	4.1%
Blocks in large tracts (MOS >1200, TRACTMOS <sup>2</sup> ≥ 2000) where Mailable addresses <sup>3</sup> ≥ 75% and predicted levels of completed mail and CATI interviews prior to CAPI subsampling > 60%	1.6%	2.0%
Other Blocks in large tracts (MOS >1200, TRACTMOS ≥ 2000)	1.7%	
All other blocks (MOS >1200, TRACTMOS < 2000) where Mailable addresses ≥ 75% and predicted levels of completed mail and CATI interviews prior to CAPI subsampling > 60%	2.1%	2.7%
All other blocks (MOS >1200, TRACTMOS < 2000)	2.3%	

<sup>1</sup>MOS = Measure of size.

<sup>2</sup>TRACTMOS = Census Tract measure of size.

<sup>3</sup>Mailable addresses: Addresses that have sufficient information to be delivered by the U.S. Postal Service (as determined by ACS).

2. Second-Stage Sample Selection – Subsampling the Unmailable and Non-Responding Addresses

All addresses determined to be unmailable are subsampled for the CAPI phase of data collection at a rate of 2-in-3. Unmailable addresses, which include Remote Alaska addresses, do not go to the CATI phase of data collection. Subsequent to CATI, all addresses for which no response has been obtained prior to CAPI are subsampled based on the expected rate of completed interviews at the tract level using the following rates.

**Table 2. Second-Stage (CAPI) Subsampling Rates for the United States and Puerto Rico**

Address and Tract Characteristics	CAPI Subsampling Rate
<b>United States</b>	
Unmailable addresses and addresses in Remote Alaska	2-in-3
Mailable addresses in tracts with predicted levels of completed mail and CATI interviews prior to CAPI subsampling between 0% and 35%	1-in-2
Mailable addresses in tracts with predicted levels of completed mail and CATI interviews prior to CAPI subsampling greater than 35% and less than 51%	2-in-5
Mailable addresses in other tracts	1-in-3
<b>Puerto Rico</b>	
Unmailable addresses	2-in-3
Mailable addresses – June through December	1-in-2
Mailable addresses – January through May	1-in-3

**ESTIMATION PROCEDURE**

The estimates that appear in this product were obtained from a ratio estimation procedure that resulted in the assignment of two sets of weights: a weight to each sample person record and a weight to each sample housing unit record. Estimates of person characteristics were based on the person weight. Estimates of family, household, and housing unit characteristics were based on the housing unit weight. For any given tabulation area, a characteristic total was estimated by summing the weights assigned to the persons, households, families or housing units possessing the characteristic in the tabulation area.

Each sample person or housing unit record was assigned exactly one weight to be used to produce estimates of all characteristics. For example, if the weight given to a sample person or housing unit had a value 40, all characteristics of that person or housing unit is tabulated with the weight of 40.

Weighting areas were formed by grouping counties of similar demographic and social characteristics using Census 2000 data. The characteristics considered in the formation included:

- Percent in poverty
- Percent renting
- Percent in rural areas
- Race, ethnicity, age, and sex distribution
- Distance between the centroids of the counties
- Core-based Statistical Area status

Each weighting area was also required to meet a threshold of 400 expected person interviews in the 2005 ACS. The stratification process then attempted to minimize the differences on the characteristics listed above between the counties within a weighting area. The process also tried to preserve as many counties that met the threshold to form their own weighting areas. In total, there were 2,006 weighting areas formed from the 3,219 counties and county equivalents including Puerto Rico.

The estimation procedure used to assign the weights was then performed independently within each of the ACS weighting areas.

1. Initial Housing Unit Weighting Factors—This process produced the following factors:

- Base Weight (BW)—This initial weight was assigned to every housing unit as the inverse of its block’s sampling rate.
- CAPI Subsampling Factor (SSF)—The weights of the CAPI cases were adjusted to reflect the results of CAPI subsampling. This factor was assigned to each record as follows:

Selected in CAPI subsampling: SSF = 2.0, 2.5, or 3.0 according to Table 2  
Not selected in CAPI subsampling: SSF = 0.0  
Not a CAPI case: SSF = 1.0

Some sample addresses were unmailable. A two-thirds sample of these were sent directly to CAPI and for these cases SSF = 1.5.

- Variation in Monthly Response by Mode (VMS)—This factor made the total weight of the Mail, CATI, and CAPI records to be tabulated in a month equal to the total base weight of all cases originally mailed for that month. For all cases, VMS was computed and assigned based on the following groups:

Strata × Month

- Noninterview Factor (NIF)—This factor adjusted the weight of all responding occupied housing units to account for both responding and nonresponding housing units. The factor was computed in two stages. The first factor, NIF1, is a ratio adjustment that was computed and assigned to occupied housings units based on the following groups:

Strata × Building Type × Tract

A second factor, NIF2, is a ratio adjustment that was computed and assigned to occupied housing units based on the following groups:

Strata × Building Type × Month

NIF was then computed by applying NIF1 and NIF2 for each occupied housing unit. Vacant housing units were assigned a value of  $NIF = 1.0$ . Nonresponding housing units were now assigned a weight of 0.0.

- Noninterview Factor—Mode (NIFM)—This factor adjusted the weight of just the responding CAPI occupied housing units to account for both CAPI respondents and all nonrespondents. This factor was computed as if NIF had not already been assigned to every occupied housing unit record. This factor was not used directly but rather as part of computing the next factor, the Mode Bias Factor.

NIFM was computed and assigned to occupied CAPI housing units based on the following groups:

$Strata \times Building\ Type \times Month$

Vacant housing units or non-CAPI (mail and CATI) housing units received a value of  $NIFM = 1.0$ .

- Mode Bias Factor (MBF)—This factor made the total weight of the housing units in the groups below the same as if NIFM had been used instead of NIF. MBF was computed and assigned to occupied housing units based on the following groups:

$Strata \times Tenure\ (Owner\ or\ renter) \times Month \times Marital\ Status\ of\ the\ Householder\ (married/widowed\ or\ single)$

Vacant housing units received a value of  $MBF = 1.0$ . MBF is applied to the weights computed through NIF.

- Housing unit Post-stratification Factor (HPF1)—This factor made the total weight of all housing units agree with the 2005 independent housing unit estimates at the weighting area level.

These independent housing unit estimates exist only for the U.S. and not Puerto Rico. Thus, all housing units in Puerto Rico received a value of  $HPF1 = 1.0$ .

2. Person Weighting Factors—Initially the person weight of each person in an occupied housing unit was the product of the weighting factors of their associated housing unit ( $BW \times \dots \times HPF1$ ). At this point everyone in the household has the same weight. These person weights were then individually adjusted based on each person's age, race, sex, and Hispanic origin as described below.

- Person Post-Stratification Factor (PPSF)—This factor was applied to individuals based on their age, race, sex and Hispanic origin in the U.S. and based on their age and sex in Puerto Rico. It adjusted the person weights so that the weighted sample counts matched independent population estimates by age, race, sex, and Hispanic origin at the weighting area level in the U.S. and matched the independent population estimates by age and sex in Puerto Rico at the weighting

area level. Because of collapsing of groups in applying this factor, only total population is assured of agreeing with the official 2005 intercensal population estimates at the weighting area level.

For U.S., this used the following groups:

Strata  $\times$  Race / Ethnicity (non-Hispanic White, non-Hispanic Black, non-Hispanic American Indian or Alaskan Native, non-Hispanic Asian, non-Hispanic Native Hawaiian or Pacific Islander, and Hispanic (any race))  $\times$  Sex  $\times$  Age Groups.

In Puerto Rico, this used only the Sex  $\times$  Age Groups.

- Rounding—The final product of all person weights ( $BW \times \dots \times HPF1 \times PPSF$ ) was rounded to an integer. Rounding was performed so that the sum of the rounded weights was within one person of the sum of the unrounded weights for any of the groups listed below:

County  
County  $\times$  Race  
County  $\times$  Race  $\times$  Hispanic Origin  
County  $\times$  Race  $\times$  Hispanic Origin  $\times$  Sex  
County  $\times$  Race  $\times$  Hispanic Origin  $\times$  Sex  $\times$  Age  
County  $\times$  Race  $\times$  Hispanic Origin  $\times$  Sex  $\times$  Age  $\times$  Tract  
County  $\times$  Race  $\times$  Hispanic Origin  $\times$  Sex  $\times$  Age  $\times$  Tract  $\times$  Block

For example, the number of White, Hispanic, Males, Age 30 estimated for a county using the rounded weights was within one of the number produced using the unrounded weights.

### 3. Final Housing Unit Weighting Factors—This process produced the following factors:

- Principal Person Factor (PPF)—This factor adjusted for differential response depending on the race, Hispanic origin, sex, and age of the principal person in the household. The principal person was defined as the female spouse of the householder. If there was no such person, then the householder was the principal person. The value of PPF for a housing unit was the PPSF of the principal person.

This adjustment was not performed for Puerto Rico because the application of this adjustment is intertwined with the ability to apply the next factor. Thus the value of  $PPF = 1.0$  was assigned to all housing units in Puerto Rico.

- Final Housing Unit Controls (HPF2)—The final product of the principal person weights ( $BW \times \dots \times HPF1 \times PPF$ ) was then assigned to the housing unit. The total number of weighted housing unit counts are then made to agree to the 2005 independent housing unit estimates at the weighting area level.



Like HPF1, this adjustment was performed only for the U.S. and not Puerto Rico since no independent housing unit estimates exist for Puerto Rico. The value of HPF2 = 1.0 was set for all housing units in Puerto Rico.

- Rounding—The final product of all housing unit weights ( $BW \times \dots \times PPF \times HPF2$ ) was rounded to an integer. Rounding was performed so that total rounded weight was within one housing unit of the total unrounded weight for any of the groups listed below:

County  
County  $\times$  Tract  
County  $\times$  Tract  $\times$  Block

## CONFIDENTIALITY OF THE DATA

The Census Bureau has modified or suppressed some data on this site to protect confidentiality. Title 13 United States Code, Section 9, prohibits the Census Bureau from publishing results in which an individual's data can be identified.

The Census Bureau's internal Disclosure Review Board sets the confidentiality rules for all data releases. A checklist approach is used to ensure that all potential risks to the confidentiality of the data are considered and addressed.

- Title 13, United States Code: Title 13 of the United States Code authorizes the Census Bureau to conduct censuses and surveys. Section 9 of the same Title requires that any information collected from the public under the authority of Title 13 be maintained as confidential. Section 214 of Title 13 and Sections 3559 and 3571 of Title 18 of the United States Code provide for the imposition of penalties of up to five years in prison and up to \$250,000 in fines for wrongful disclosure of confidential census information.
- Disclosure Limitation: Disclosure limitation is the process for protecting the confidentiality of data. A disclosure of data occurs when someone can use published statistical information to identify an individual that has provided information under a pledge of confidentiality. For data tabulations the Census Bureau uses disclosure limitation procedures to modify or remove the characteristics that put confidential information at risk for disclosure. Although it may appear that a table shows information about a specific individual, the Census Bureau has taken steps to disguise or suppress the original data while making sure the results are still useful. The techniques used by the Census Bureau to protect confidentiality in tabulations vary, depending on the type of data.
- Data Swapping: Data swapping is a method of disclosure limitation designed to protect confidentiality in tables of frequency data (the number or percent of the population with

certain characteristics). Data swapping is done by editing the source data or exchanging records for a sample of cases when creating a table. A sample of households is selected and matched on a set of selected key variables with households in neighboring geographic areas that have similar characteristics (such as the same number of adults and same number of children). Because the swap often occurs within a neighboring area, there is no effect on the marginal totals for the area or for totals that include data from multiple areas. Because of data swapping, users should not assume that tables with cells having a value of one or two reveal information about specific individuals. Data swapping procedures were first used in the 1990 Census, and were used again in Census 2000.

## ERRORS IN THE DATA

- **Sampling Error** — The data in the ACS products are estimates of the actual figures that would have been obtained by interviewing the entire population using the same methodology. The estimates from the chosen sample also differ from other samples of housing units and persons within those housing units. Sampling error in data arises due to the use of probability sampling, which is necessary to ensure the integrity and representativeness of sample survey results. The implementation of statistical sampling procedures provides the basis for the statistical analysis of sample data.
- **Nonsampling Error** — In addition to sampling error, data users should realize that other types of errors may be introduced during any of the various complex operations used to collect and process survey data. For example, operations such as data entry from questionnaires and editing may introduce error into the estimates. These and other sources of error contribute to the nonsampling error component of the total error of survey estimates. Nonsampling errors may affect the data in two ways. Errors that are introduced randomly increase the variability of the data. Systematic errors which are consistent in one direction introduce bias into the results of a sample survey. The Census Bureau protects against the effect of systematic errors on survey estimates by conducting extensive research and evaluation programs on sampling techniques, questionnaire design, and data collection and processing procedures. In addition, an important goal of the ACS is to minimize the amount of nonsampling error introduced through nonresponse for sample housing units. One way of accomplishing this is by following up on mail nonrespondents during the CATI and CAPI phases.

## MEASURES OF SAMPLING ERROR

Sampling error is the difference between an estimate based on a sample and the corresponding value that would be obtained if the estimate were based on the entire population (as from a census). Note that sample-based estimates will vary depending on the particular sample selected

from the population. Measures of the magnitude of sampling error reflect the variation in the estimates over all possible samples that could have been selected from the population using the same sampling methodology.

Estimates of the magnitude of sampling errors – in the form of margins of error – are provided with all published ACS data. The Census Bureau recommends that data users incorporate this information into their analyses, as sampling error in survey estimates could impact the conclusions drawn from the results.

### Confidence Intervals and Margins of Error

Confidence Intervals – A sample estimate and its estimated standard error may be used to construct confidence intervals about the estimate. These intervals are ranges that will contain the average value of the estimated characteristic that results over all possible samples, with a known probability.

For example, if all possible samples that could result under the ACS sample design were independently selected and surveyed under the same conditions, and if the estimate and its estimated standard error were calculated for each of these samples, then:

1. Approximately 68 percent of the intervals from one estimated standard error below the estimate to one estimated standard error above the estimate would contain the average result from all possible samples;
2. Approximately 90 percent of the intervals from 1.65 times the estimated standard error below the estimate to 1.65 times the estimated standard error above the estimate would contain the average result from all possible samples.
3. Approximately 95 percent of the intervals from two estimated standard errors below the estimate to two estimated standard errors above the estimate would contain the average result from all possible samples.

The intervals are referred to as 68 percent, 90 percent, and 95 percent confidence intervals, respectively.

Margin of Error – Instead of providing the upper and lower confidence bounds in published ACS tables, the margin of error is provided instead. The margin of error is the difference between an estimate and its upper or lower confidence bound. Both the confidence bounds and the standard error can easily be computed from the margin of error. All ACS published margins of error are based on a 90 percent confidence level.

$$\text{Standard Error} = \text{Margin of Error} / 1.65$$

$$\text{Lower Confidence Bound} = \text{Estimate} - \text{Margin of Error}$$

Upper Confidence Bound = Estimate + Margin of Error

When constructing confidence bounds from the margin of error, the user should be aware of any “natural” limits on the bounds. For example, if a population estimate is near zero, the calculated value of the lower confidence bound may be negative. However, a negative number of people does not make sense, so the lower confidence bound should be reported as zero instead. However, for other estimates such as income, negative values do make sense. The context and meaning of the estimate must be kept in mind when creating these bounds. Another of these natural limits would be 100% for the upper bound of a percent estimate.

If the margin of error is displayed as ‘\*\*\*\*\*’ (five asterisks), the estimate has been controlled to be equal to a fixed value and so has no sampling error. When using any of the formulas in the following section, use a standard error of zero for these controlled estimates.

Limitations –The user should be careful when computing and interpreting confidence intervals.

- The estimated standard errors included in these data products do not include portions of the variability due to nonsampling error that may be present in the data. In particular, the standard errors do not reflect the effect of correlated errors introduced by interviewers, coders, or other field or processing personnel. Nor do they reflect the error from imputed values due to missing responses. Thus, the standard errors calculated represent a lower bound of the total error. As a result, confidence intervals formed using these estimated standard errors may not meet the stated levels of confidence (i.e., 68, 90, or 95 percent). Thus, some care must be exercised in the interpretation of the data in this data product based on the estimated standard errors.
- Zero or small estimates; very large estimates — The value of almost all ACS characteristics is greater than or equal to zero by definition. For zero or small estimates, use of the method given previously for calculating confidence intervals relies on large sample theory, and may result in negative values which for most characteristics are not admissible. In this case the lower limit of the confidence interval is set to zero by default. A similar caution holds for estimates of totals close to a control total or estimated proportions near one, where the upper limit of the confidence interval is set to its largest admissible value. In these situations the level of confidence of the adjusted range of values is less than the prescribed confidence level.

## CALCULATION OF STANDARD ERRORS

Direct estimates of the standard errors were calculated for all estimates reported in this product. The standard errors, in most cases, are calculated using a replicate-based methodology that takes into account the sample design and estimation procedures. Exceptions include:

1. The estimate of the number or proportion of people, households, families, or housing units in a geographic area with a specific characteristic is zero. A special procedure is used to estimate the standard error.
2. There are no sample observations available to compute an estimate of a median, a proportion, or some other ratio, or an estimate of its standard error. The estimate is represented in the tables by “-” and the margin of error by “\*\*” (two asterisks).
3. Only a small number of identical values are reported and used to calculate a median, aggregate, mean, or per capita amount. In this case, there are too few sample observations to compute a stable estimate of the standard error. The margin of error is represented in the tables by “\*” (one asterisk).
4. The estimate of a median falls in the lower open-ended interval or upper open-ended interval of a distribution. If the median occurs in the lowest interval, then a “-” follows the estimate, and if the median occurs in the upper interval, then a “+” follows the estimate. In both cases the margin of error is represented in the tables by “\*\*\*” (three asterisks).

Sums and Differences of Direct Standard Errors — The standard errors estimated from these tables are for individual estimates. Additional calculations are required to estimate the standard errors for sums of and differences between two sample estimates. The estimate of the standard error of a sum or difference is approximately the square root of the sum of the two individual standard errors squared; that is, for standard errors  $SE(\hat{X})$  and  $SE(\hat{Y})$  of estimates  $\hat{X}$  and  $\hat{Y}$ :

$$SE(\hat{X} + \hat{Y}) = SE(\hat{X} - \hat{Y}) = \sqrt{[SE(\hat{X})]^2 + [SE(\hat{Y})]^2}$$

This method, however, will underestimate (overestimate) the standard error if the two items in a sum are highly positively (negatively) correlated or if the two items in a difference are highly negatively (positively) correlated.

Ratios — The statistic of interest may be the ratio of two estimates. First is the case where the numerator *is not* a subset of the denominator. The standard error of this ratio between two sample estimates is approximated as:

$$SE\left(\frac{\hat{X}}{\hat{Y}}\right) = \frac{1}{\hat{Y}} \sqrt{[SE(\hat{X})]^2 + \frac{\hat{X}^2}{\hat{Y}^2} [SE(\hat{Y})]^2}$$

Proportions/percentages – For a proportion (or percent), a ratio where the numerator *is* a subset of the denominator, a slightly different estimator is used. Note the difference between the formulas for the standard error for proportions (below) and ratios (above) - the plus sign in the previous

formula has been replaced with a minus sign. If the value under the square root sign is negative, use the ratio standard error formula above, instead. If  $\hat{P} = \hat{X} / \hat{Y}$ , then

$$SE(\hat{P}) = \frac{1}{\hat{Y}} \sqrt{[SE(\hat{X})]^2 - \frac{\hat{X}^2}{\hat{Y}^2} [SE(\hat{Y})]^2}$$

If  $\hat{Q} = 100\% \times \hat{P}$  (P is the proportion and Q is its corresponding percent), then  $SE(\hat{Q}) = 100\% \times SE(\hat{P})$ .

Products – For a product of two estimates - for example if you want to estimate a proportion's numerator by multiplying the proportion by its denominator - the standard error can be approximated as

$$SE(\hat{X} \times \hat{Y}) = \sqrt{\hat{X}^2 \times [SE(\hat{Y})]^2 + \hat{Y}^2 \times [SE(\hat{X})]^2}$$

Significant differences – Users may conduct a statistical test to see if the difference between an ACS estimate and any other chosen estimates is statistically significant at a given confidence level. “Statistically significant” means that the difference is not likely due to random chance alone. With the two estimates (Est<sub>1</sub> and Est<sub>2</sub>) and their respective standard errors (SE<sub>1</sub> and SE<sub>2</sub>), calculate

$$Z = \frac{Est_1 - Est_2}{\sqrt{(SE_1)^2 + (SE_2)^2}}$$

If  $Z > 1.65$  or  $Z < -1.65$ , then the difference can be said to be statistically significant at the 90% confidence level. Any estimate can be compared to an ACS estimate using this method, including other ACS estimates from the current year, the ACS estimate for the same characteristic and geographic area but from a previous year, Census 2000 100% counts and long form estimates, estimates from other Census Bureau surveys, and estimates from other sources. Not all estimates have sampling error — Census 2000 100% counts do not, for example, although Census 2000 long form estimates do — but they should be used if they exist to give the most accurate result of the test.

Users are also cautioned to *not* rely on looking at whether confidence intervals for two estimates overlap to determine statistical significance, because there are circumstances where that method will not give the correct test result. The Z calculation above is recommended in all cases.

All statistical testing in ACS data products is based on the 90% confidence level. Users should understand that all testing was done using *unrounded* estimates and standard errors, and it may not be possible to replicate test results using the rounded estimates and margins of error as published.

## EXAMPLES OF STANDARD ERROR CALCULATIONS

We will present some examples based on the real data to demonstrate the use of the formulas.

### Example 1 - Calculating the Standard Error from the Confidence Interval

The estimated number of males, never married is 34,171,130 from summary table B12001 for the United States for 2005. The margin of error is 81,645.

$$\text{Standard Error} = \text{Margin of Error} / 1.65$$

Calculating the standard error using the margin of error, we have:

$$\text{SE}(34,171,130) = 81,645 / 1.65 = 49,482.$$

### Example 2 - Calculating the Standard Error of a Sum

We are interested in the number of people who have never been married. From Example 1, we know the number of males, never married is 34,171,130. From summary table B12001 we have the number of females, never married is 29,943,646 with a margin of error of 74,944. So, the estimated number of people who have never been married is  $34,171,130 + 29,943,646 = 64,114,776$ . To calculate the standard error of this sum, we need the standard errors of the two estimates in the sum. We have the standard error for the number of males never married from example 1 as 49,482. The standard error for the number of females never married is calculated using the margin of error:

$$\text{SE}(29,943,646) = 74,944 / 1.65 = 45,421.$$

So using the formula for the standard error of a sum or difference we have:

$$\text{SE}(64,114,776) = \sqrt{49,482^2 + 45,421^2} = 67,168$$

Caution: This method, however, will underestimate (overestimate) the standard error if the two items in a sum are highly positively (negatively) correlated or if the two items in a difference are highly negatively (positively) correlated.

To calculate the lower and upper bounds of the 90 percent confidence interval around 64,114,776 using the standard error, simply multiply 67,168 by 1.65, then add and subtract the product from 64,114,776. Thus the 90 percent confidence interval for this estimate is  $[64,114,776 - 1.65(67,168)]$  to  $[64,114,776 + 1.65(67,168)]$  or 64,003,949 to 64,225,603.

### Example 3 - Calculating the Standard Error of a Percent

We are interested in the percentage of females who have never been married to the number of people who have never been married. The number of females, never married is 29,943,646 and the number of people who have never been married is 64,114,776. To calculate the standard error of this sum, we need the standard errors of the two estimates in the sum. We have the standard error for the number of females never married from example 2 as 45,421 and the standard error for the number of people never married calculated from example 2 as 67,168.

The estimate is  $(29,943,646 / 64,114,776) * 100\% = 46.7\%$

So, using the formula for the standard error of a proportion or percent, we have:

$$SE(46.7\%) = 100\% * \left( \frac{1}{64,114,776} \sqrt{45,521^2 - 0.467^2 \times 67,168^2} \right) = 0.05\%$$

To calculate the lower and upper bounds of the 90 percent confidence interval around 46.7 using the standard error, simply multiply 0.05 by 1.65, then add and subtract the product from 46.7. Thus the 90 percent confidence interval for this estimate is  $[46.7 - 1.65(0.05)]$  to  $[46.7 + 1.65(0.05)]$ , or 46.6% to 46.8%.

### CONTROL OF NONSAMPLING ERROR

As mentioned earlier, sample data are subject to nonsampling error. This component of error could introduce serious bias into the data, and the total error could increase dramatically over that which would result purely from sampling. While it is impossible to completely eliminate nonsampling error from a survey operation, the Census Bureau attempts to control the sources of such error during the collection and processing operations. Described below are the primary sources of nonsampling error and the programs instituted for control of this error. The success of these programs, however, is contingent upon how well the instructions were carried out during the survey.

- Undercoverage — It is possible for some sample housing units or persons to be missed entirely by the survey. The undercoverage of persons and housing units can introduce biases into the data. A major way to avoid undercoverage in a survey is to ensure that its sampling frame, for ACS an address list in each state, is as complete and accurate as possible.

The source of addresses was the MAF. The MAF is created by combining the Delivery Sequence File of the United States Postal Service, and the address list for Census 2000. An attempt is made to assign all appropriate geographic codes to each MAF address via an automated procedure using the Census Bureau TIGER files. A manual coding



operation based in the appropriate regional offices is attempted for addresses which could not be automatically coded. The MAF was used as the source of addresses for selecting sample housing units and mailing questionnaires. TIGER produced the location maps for CAPI assignments.

In the CATI and CAPI nonresponse follow-up phases, efforts were made to minimize the chances that housing units that were not part of the sample were interviewed in place of units in sample by mistake. If a CATI interviewer called a mail nonresponse case and was not able to reach the exact address, no interview was conducted and the case was eligible for CAPI. During CAPI follow-up, the interviewer had to locate the exact address for each sample housing unit. If the interviewer could not locate the exact sample unit in a multi-unit structure, or found a different number of units than expected, the interviewers were instructed to list the units in the building and follow a specific procedure to select a replacement sample unit.

- Respondent and Interviewer Error — The person completing the questionnaire or responding to the questions posed by an interviewer could serve as a source of error, although the questions were cognitively tested for phrasing, and detailed instructions for completing the questionnaire were provided to each household.
  - Interviewer monitoring — The interviewer may misinterpret or otherwise incorrectly enter information given by a respondent; may fail to collect some of the information for a person or household; or may collect data for households that were not designated as part of the sample. To control these problems, the work of interviewers was monitored carefully. Field staff were prepared for their tasks by using specially developed training packages that included hands-on experience in using survey materials. A sample of the households interviewed by CAPI interviewers was reinterviewed to control for the possibility that interviewers may have fabricated data.
  - Item Nonresponse — Nonresponse to particular questions on the survey questionnaire and instrument allows for the introduction of bias into the data, since the characteristics of the nonrespondents have not been observed and may differ from those reported by respondents. As a result, any imputation procedure using respondent data may not completely reflect this difference either at the elemental level (individual person or housing unit) or on average.

Some protection against the introduction of large biases is afforded by minimizing nonresponse. In the ACS, item nonresponse for the CATI and CAPI operations was minimized by the requirement that the automated instrument receive a response to each question before the next one could be asked. Questionnaires returned by mail were edited for completeness and acceptability. They were reviewed by computer for content omissions and population coverage. If necessary, a telephone follow-up was

made to obtain missing information. Potential coverage errors were included in this follow-up.

- **Processing Error** — The many phases involved in processing the survey data represent potential sources for the introduction of nonsampling error. The processing of the survey questionnaires includes the keying of data from completed questionnaires, automated clerical review, and follow-up by telephone, the manual coding of write-in responses, and the electronic data processing. The various field, coding and computer operations undergo a number of quality control checks to insure their accurate application.
- **Content Editing** — After data collection was completed, any remaining incomplete or inconsistent information was imputed during the final content edit of the collected data. Imputations, or computer assignments of acceptable codes in place of unacceptable entries or blanks, were needed most often when an entry for a given item was missing or when the information reported for a person or housing unit on that item was inconsistent with other information for that same person or housing unit. As in other surveys and previous censuses, the general procedure for changing unacceptable entries was to allocate an entry for a person or housing unit that was consistent with entries for persons or housing units with similar characteristics. Imputing acceptable values in place of blanks or unacceptable entries enhances the usefulness of the data.